

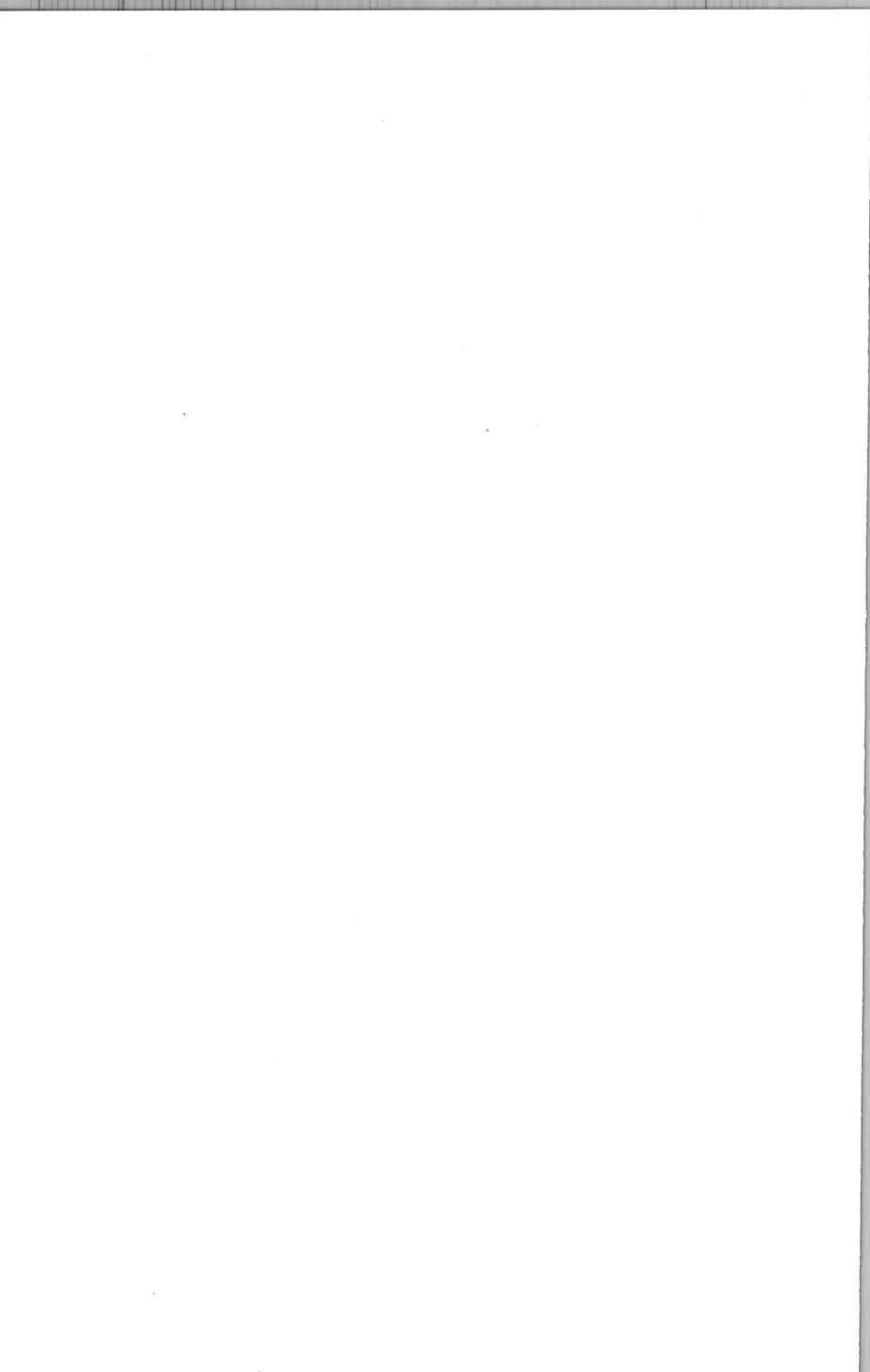
DEOS Progress Letter

Edited by Roland Klees

no 98.2

DEOS





729222

9

DEOS Progress Letter

98.2

Bibliotheek TU Delft



C 3036913

2393
245
8

Reprint of internal DEOS Progress Letter 97.1

DEOS Progress Letter

98.2

Edited by Roland Klees



Published and distributed by:

Delft University Press
P.O. Box 98
2600 MG Delft
The Netherlands
Telephone: +31 15 2783254
Telefax: + 31 15 2781661
E-mail: DUP@DUP.TUdelft.NL

ISBN 90-407-1820-2

Copyright 1998 by DEOS

All rights reserved. No part of the material protected by this copyright notice may be reproduced or utilized in any form or by any means, electronic or mechanical, including photocopying, recording or by any information storage and retrieval system, without written permission from the publisher: Delft University Press.

Printed in The Netherlands

CONTENTS

Integral Equation Formulations for Geodetic Mixed Boundary Value Problems	1
Roland Klees, Stefan Ritter, Rüdiger Lehmann	
On Singular Surface Integrals in Physical Geodesy	15
Roland Klees	
The Fiction of the Geoid	29
Roland Klees, Martin van Gelderen	
Natural Gas Extraction and its Induced Gravity Change	45
Martin van Gelderen, Roger Haagmans, Mirjam Bilker	
Error Propagation for Satellite Gradiometry	57
Martin van Gelderen	
Quality Differences between Tikhonov Regularization and Generalized Biased Estimation in Gradiometric Analysis	69
Johannes Bouman, Radboud Koop	
Quality Assessment of Geopotential Models by Means of Redundancy Composition?	81
Johannes Bouman	
Overview of Tide Gauge Systems and Averaging Techniques	91
Kyra van Onselen	
Research Plan and Progress Report of the TOPEX/POSEIDON Extended Mission	111
Marc Naeije, Roger Haagmans, Ernst Schrama, Karel Wakker, Remko Scharroo	

1
2
3
4
5
6
7
8
9
10
11
12
13
14
15
16
17
18
19
20
21
22
23
24
25
26
27
28
29
30
31
32
33
34
35
36
37
38
39
40
41
42
43
44
45
46
47
48
49
50
51
52
53
54
55
56
57
58
59
60
61
62
63
64
65
66
67
68
69
70
71
72
73
74
75
76
77
78
79
80
81
82
83
84
85
86
87
88
89
90
91
92
93
94
95
96
97
98
99
100

Integral Equation Formulations for Geodetic Mixed Boundary Value Problems

Roland Klees, Stefan Ritter¹, Rüdiger Lehmann²

¹ Mathematical Institute, University of Karlsruhe, Germany

² Geodetic Institute, University of Karlsruhe, Germany

Abstract

We consider mixed boundary value problems in Physical Geodesy and study possibilities in order to transform them into a system of integral equations over the boundary of the domain. The system of integral equations can be solved numerically, by, e.g. boundary element methods, provided that (a) the mixed boundary value problem is uniquely solvable, (b) the system of integral equations is equivalent to the mixed boundary value problem, and (c) the matrix of integral operators is strongly elliptic. We introduce a method, first proposed by Stephan, which allows to derive integral equation formulations for all mixed boundary value problems relevant to geodetic applications. Moreover, the analysis of Stephan for the mixed Dirichlet-Neumann problem may be generalized to the geodetic mixed boundary value problems, as well.

1 Introduction

The objective of the paper is to study mixed boundary value problems (MBVPs) of type

$$\begin{aligned}
 \Delta u &= 0 && \text{in } D^c \\
 B_o u &= g_o && \text{on } S_o \\
 B_c u &= g_c && \text{on } S_c \\
 u &= O(|x|^{-1}), && |x| \rightarrow \infty.
 \end{aligned}
 \tag{1}$$

D is a bounded domain in \mathbb{R}^3 with sufficiently smooth boundary $S = S_o \cup S_c$, with $S_o \cap S_c = \emptyset$, and D^c its complement in \mathbb{R}^3 , i.e. $D^c = \mathbb{R}^3 \setminus \bar{D}$. The curve $\bar{S}_o \cap \bar{S}_c$ is assumed to be smooth and simply closed. B_o and B_c are first-order differential operators, and g_o and g_c are the given boundary data. In geodetic applications, S_o can be identified with the surface of the oceans and S_c with the continents. Depending of the choice of B_o and B_c different mixed boundary value problems can be formulated. In geodesy, the most relevant (linearized) MBVPs are summarized in Table 1. Depending on the level of approximation, additional MBVPs can be derived from

Table 1. Linearized geodetic mixed boundary value problems

name	B_o	B_c	type
altimetry-gravimetry I	I	$I - D_\tau$	Dirichlet-Poincaré
altimetry-gravimetry II	D_τ	$I - D_\tau$	oblique-Poincaré
fixed altimetry-gravimetry	I	D_τ	Dirichlet-oblique

the three basic problems listed in Table 1. For instance, in spherical approximation and constant radius approximation, the oblique boundary operator D_τ becomes the Neumann operator, and the Poincaré boundary operator $I - D_\tau$ becomes the Robin operator.

Existence and uniqueness of various linearized geodetic MBVPs have been studied, mostly in the context of the spherical and constant radius approximation, see, e.g., Arnold (1981); Sjöberg (1982); Holota (1982); Svensson (1983); Sacerdote and Sansò (1983a, 1983b); Holota (1983a, 1983b); Arnold (1984); Sacerdote and Sansò (1987); Svensson (1988); Sansò (1993); Keller (1996).

Numerical aspects of geodetic MBVPs have been studied by, e.g., Sjöberg (1982); Bjerhammar (1983); Sansò and Stock (1985); Hofmann-Wellenhof (1985); Mainville (1986); Mayer (1997). In the context of integral equation formulations, the references Sansò and Stock (1985) and Mayer (1997) are of interest. In Sansò and Stock (1985) an integral equation formulation of the linearized altimetry-gravimetry II MBVP in spherical approximation has been used and applied to a local area (see Section 4). The transformation of the MBVP into an integral equation is based on the explicit solution of the Neumann problem for a spherical boundary surface S , and cannot be applied to MBVPs with non-spherical surfaces and/or other types of boundary data. Mayer (1997) proposes a completely new solution strategy for the nonlinear altimetry-gravimetry II MBVP, which assumes a global coverage with gravity values and, in addition, the potential to be given on the free continental part of the boundary. Firstly, a hypersingular integral equation for the linearized fixed gravimetric BVP is solved, due to the global coverage with gravity values. Then, the remaining Dirichlet boundary condition over the free continental part yields a nonlinear operator equation, which has to be solved for the unknown continental geometry. The solution has to be improved iteratively (see Section 5).

Stephan (1987) studied the Dirichlet-Neumann MBVP on closed surfaces in \mathbb{R}^3 based on an equivalent formulation of the MBVP as a system of two integral equations. His method is general enough to derive integral equations for all relevant MBVPs in geodesy. Moreover, his procedure to prove the existence and uniqueness of the system of integral equations, and the equivalence of the MBVP with the system of integral equations, may be applied to geodetic MBVP, as well. Therefore we first want to introduce his method and the main lines of the analysis; then we want to show how integral equations for geodetic MBVPs can be derived analogously. Finally, we will briefly discuss the methods of Sansò and Stock (1985) and Mayer (1997) since they also make use of integral equations in order to solve geodetic MBVPs.

2 The method of Stephan

Stephan (1987) discusses the solution of the Dirichlet-Neumann problem in \mathbb{R}^3 :

$$\begin{aligned} \Delta u &= 0 && \text{in } D^c \\ u &= g_1 && \text{on } S_1 \\ D_n u &= g_2 && \text{on } S_2 \\ u &= O(|x|^{-1}), && |x| \rightarrow \infty. \end{aligned} \quad (2)$$

n is the unit normal vector to S pointing into D^c . S is assumed to be sufficiently smooth. An extension to polyhedral domains is presented in von Petersdorff and Stephan (1990).

Existence and uniqueness of the weak solution u of the Dirichlet-Neumann MBVP (2) is proved by use of (a) the uniqueness of the weak solution, (b) the equivalence of the MBVP to the system of integral equations, (c) the existence and uniqueness of the solution of the system of integral equations, and (d) the solution of the integral equation by inserting into the representation formula. The *weak solution* is defined by Green's first identity: Let $u \in H_{loc}^1(D^c)$ and $v \in H^1(D^c)$ with bounded support. Then

$$\int_{D^c} \nabla u \nabla v \, dD = - \int_S \gamma(D_n u) \gamma v \, dS, \quad (3)$$

where γ denotes the restriction to S . This holds if the trace $\gamma D_n u$ is at least in $H^{-1/2}(S)$. The space U we look for the weak solution is defined by $U := \{u \in H_{loc}^1(D^c) : \Delta u = 0 \text{ in } D^c, u = O(|x|^{-1}), |x| \rightarrow \infty\}$.

The *uniqueness of the weak solution* can easily be proved by means of Green's first identity applied to $\Omega := B \cap D^c$, where B denotes a sufficiently large ball with radius R that encloses \bar{D} . Let $u \in U$ with $\gamma_1 u = 0$ and $\gamma_2 D_n u = 0$, where γ_i is the restriction to S_i , $i = \{1, 2\}$. Then

$$\int_{\partial B} u D_n u \, d\partial B = - \int_{\Omega} |\nabla u|^2 \, d\Omega. \quad (4)$$

The left-hand side of (4) tends to zero as $R \rightarrow \infty$. This implies $|\nabla u| = 0$, thus $u = \text{const.}$ in Ω . Since $u = O(|x|^{-1})$, it follows that $u = 0$.

In order to transform the MBVP (2) into an integral equation we need a representation of the weak solution u of the MBVP in terms of boundary potentials. This can be done in different ways, e.g., by using a representation of u as single layer potential, double layer potential or a linear combination of both. Here, we make use of another possibility, namely of Green's third identity: For $u \in U$, and the *Cauchy-data*

$$\begin{pmatrix} \mu \\ \nu \end{pmatrix} := \begin{pmatrix} \gamma u \\ \gamma D_n u \end{pmatrix} \in H^{1/2}(S) \times H^{-1/2}(S), \quad (5)$$

it holds

$$u(x) = \int_S \mu(y) D_{n(y)} s(x-y) \, dS(y) - \int_S s(x-y) \nu(y) \, dS(y), \quad x \in D^c, \quad (6)$$

where $s(x - y)$ is the fundamental solution of the Laplace equation in \mathbb{R}^3 , i.e., $s(x - y) = (4\pi |y - x|)^{-1}$. The *Calderon-projector*

$$P := \left(\frac{1}{2}I + A \right) \quad \text{with} \quad A := \begin{pmatrix} K & -V \\ D & -K' \end{pmatrix} \quad (7)$$

projects $H^{1/2}(S) \times H^{1/2}(S)$ on the Cauchy-data of weak solutions in U , see Stephan (1987). This projector might be understood as generalization of the well-known limit-relations for the single layer potential and the double layer potential (e.g., Miranda (1970)). Using $P \begin{pmatrix} \mu \\ \nu \end{pmatrix}$ for the Cauchy-data of the weak solution, we obtain the following system of integral equations on S :

$$\frac{1}{2} \begin{pmatrix} \mu \\ \nu \end{pmatrix} = \begin{pmatrix} K & -V \\ D & -K' \end{pmatrix} \begin{pmatrix} \mu \\ \nu \end{pmatrix}. \quad (8)$$

The boundary integral operators V , K , K' , and D are defined by ($x \in S$):

$$\begin{aligned} (V\chi)(x) &= \int_S s(x - y) \chi(y) dS(y), \\ (K\chi)(x) &= \int_S D_{n(y)} s(x - y) \chi(y) dS(y), \\ (K'\chi)(x) &= \int_S D_{n(x)} s(x - y) \chi(y) dS(y), \\ (D\chi)(x) &= \int_S D_{n(x)} D_{n(y)} s(x - y) \chi(y) dS(y). \end{aligned} \quad (9)$$

The system (8) together with the boundary conditions in (2) provide more equations than unknowns; depending on how they are combined, we can derive a system of first order integral equations, of second order integral equations, or a mixed system of integral equations. For instance, when restricting the first equation in (8) to S_1 and the second equation to S_2 , we obtain

$$\begin{aligned} \text{on } S_1: \quad \frac{1}{2}g_1 &= K_{11}g_1 + K_{21}\mu - V_{11}\nu - V_{21}g_2, \\ \text{on } S_2: \quad \frac{1}{2}g_2 &= D_{12}g_1 + D_{22}\mu - K'_{12}\nu - K'_{22}g_2, \end{aligned} \quad (10)$$

or, in matrix form,

$$\begin{pmatrix} K_{21} & -V_{11} \\ D_{22} & -K'_{12} \end{pmatrix} \begin{pmatrix} \mu \\ \nu \end{pmatrix} = \begin{pmatrix} \frac{1}{2}I - K_{11} & V_{21} \\ -D_{12} & \frac{1}{2}I + K'_{22} \end{pmatrix} \begin{pmatrix} g_1 \\ g_2 \end{pmatrix}. \quad (11)$$

The subscript ik means integration over S_i and evaluation on S_k , e.g., if

$$(K\chi)(x) = \int_S D_{n(y)} s(x - y) \chi(y) dS(y), \quad x \in S, \quad (12)$$

the operator K_{ik} is defined by

$$(K_{ik}\chi)(x) = \int_{S_i} D_{n(y)} s(x - y) \chi(y) dS(y), \quad x \in S_k. \quad (13)$$

Equation (11) defines a system of first order integral equations for the Cauchy-data $\begin{pmatrix} \mu \\ \nu \end{pmatrix}$. Alternatively, we may restrict the first equation in (8) to S_2 and the second one to S_1 ; then we obtain

$$\begin{aligned} \text{on } S_1: \quad \frac{1}{2}\nu &= D_{11}g_1 + D_{21}\mu - K'_{11}\nu - K'_{21}g_2, \\ \text{on } S_2: \quad \frac{1}{2}\mu &= K_{12}g_1 + K_{22}\mu - V_{12}\nu - V_{22}g_2, \end{aligned} \quad (14)$$

i.e., a system of second kind integral equations

$$\begin{pmatrix} \frac{1}{2}I - K_{22} & V_{12} \\ -D_{21} & \frac{1}{2}I + K'_{11} \end{pmatrix} \begin{pmatrix} \mu \\ \nu \end{pmatrix} = \begin{pmatrix} K_{12} & -V_{22} \\ D_{11} & -K'_{21} \end{pmatrix} \begin{pmatrix} g_1 \\ g_2 \end{pmatrix}. \quad (15)$$

Analogously, mixed systems can be obtained by taking only one of the two equations in (8) and restricting first to S_1 and then to S_2 .

The *solvability* of the systems of integral equations is shown in several steps. We will omit the details and will only point out the main lines, whereby we limit to system (11). For the details, see Stephan (1987).

Firstly, the mapping properties of the involved integral operators are determined. The operators V , K , and K' have kernels of order $O(|y-x|^{-1})$ as $(y-x) \rightarrow 0$, hence they are weakly singular integral operators on S . D has a kernel of order $O(|y-x|^{-3})$ as $(y-x) \rightarrow 0$, i.e., it is a hypersingular integral operator on S . Moreover, V , K , K' , and D are continuous mappings in suitable Sobolev spaces, i.e., they define pseudodifferential operators of integer order on S . V , K , and K' have order -1 , and D has order $+1$. Although the mappings

$$\begin{aligned} V_{ik} : \tilde{H}^s(S_i) &\rightarrow H^{s+1}(S_k), \\ K_{ik} : \tilde{H}^s(S_i) &\rightarrow H^{s+1}(S_k), \\ K'_{ik} : \tilde{H}^s(S_i) &\rightarrow H^{s+1}(S_k), \\ D_{ik} : \tilde{H}^{s+1}(S_i) &\rightarrow H^s(S_k), \end{aligned} \quad (16)$$

act only on pieces of S , it can be shown, using standard arguments from the theory of pseudodifferential operators, that they are continuous for some real s , depending on the smoothness of S . Here, $u \in \tilde{H}^s(S_i) = \{u \in H^s(S) : \text{supp } u \subset \bar{S}_i\}$.

Secondly, the system (11) is rewritten in order to make use of the mapping properties (16): if $\tilde{g}_1 \in H^{1/2}(S)$ and $\tilde{g}_2 \in H^{-1/2}(S)$ denote arbitrary extensions of the boundary data, the unknown Cauchy-data $\begin{pmatrix} \mu \\ \nu \end{pmatrix}$ admit the form

$$\begin{pmatrix} \mu \\ \nu \end{pmatrix} = \begin{pmatrix} \mu_0 \\ \nu_0 \end{pmatrix} + \begin{pmatrix} \tilde{g}_1 \\ \tilde{g}_2 \end{pmatrix}, \quad (17)$$

with $\mu_0 \in \tilde{H}^{1/2}(S_2)$ and $\nu_0 \in \tilde{H}^{-1/2}(S_1)$ and $\gamma_1\mu_0 = 0$ and $\gamma_2\nu_0 = 0$. Then, (11) can be written as

$$\begin{pmatrix} K_{21} & -V_{11} \\ D_{22} & -K'_{12} \end{pmatrix} \begin{pmatrix} \mu_0 \\ \nu_0 \end{pmatrix} = \begin{pmatrix} \frac{1}{2}I - K_{S1} & V_{S1} \\ -D_{S2} & \frac{1}{2}I + K'_{S2} \end{pmatrix} \begin{pmatrix} \tilde{g}_1 \\ \tilde{g}_2 \end{pmatrix}, \quad (18)$$

where K_{S_1} means integration over S and evaluation on S_1 etc.

Thirdly, the mapping properties of the involved matrix operators

$$A := \begin{pmatrix} K_{21} & -V_{11} \\ D_{22} & -K'_{12} \end{pmatrix}, \quad B := \begin{pmatrix} \frac{1}{2}I - K_{S_1} & V_{S_1} \\ -D_{S_2} & \frac{1}{2}I + K'_{S_2} \end{pmatrix} \quad (19)$$

are investigated. From (16) it follows that the matrix operator A is continuous for some real s depending on the smoothness of S as mapping $\tilde{H}^s(S_2) \times \tilde{H}^s(S_1) \rightarrow \tilde{H}^s(S_1) \times \tilde{H}^{s-1}(S_2)$. The continuity of the extensions \tilde{g}_i , $i = \{1, 2\}$ in $H^s(S)$ for $g_i \in H^s(S_i)$ together with the mapping properties (16), provide the continuity of B as mapping $\tilde{H}^s(S) \times \tilde{H}^{s-1}(S) \rightarrow \tilde{H}^s(S_1) \times \tilde{H}^{s-1}(S_2)$, for some real s , depending on the smoothness of S .

Fourthly, the *uniqueness* of (18) is shown. We omit the details and refer to Stephan (1987). Moreover, since the matrix operator A is strongly elliptic, i.e., it satisfies some coerciveness inequalities in appropriate Sobolev spaces, it differs by a compact perturbation from a positive definite operator. Hence, A is a Fredholm operator of index zero. For Fredholm operators of index zero it is known that injectivity implies surjectivity, thus A is bijective.

Finally, the *equivalence* of the original MBVP (2) with the system of integral equations (11) is shown, i.e., $\mu_0 = \gamma_2 u - \gamma_2 \tilde{g}_1$, $\nu_0 = \gamma_1 D_n u - \gamma_1 \tilde{g}_2$, and, conversely, u in D^c is given by

$$u(x) = \int_S \tilde{\mu}(y) D_{n(y)} s(x-y) dS(y) - \int_S s(x-y) \tilde{\nu}(y) dS(y), \quad x \in D^c, \quad (20)$$

with

$$\tilde{\mu} = \begin{cases} \mu_0 + \tilde{g}_1 & \text{on } S_2 \\ g_1 & \text{on } S_1 \end{cases}, \quad \tilde{\nu} = \begin{cases} \nu_0 + \tilde{g}_2 & \text{on } S_1 \\ g_2 & \text{on } S_2 \end{cases}, \quad (21)$$

and extensions \tilde{g}_i , $i = \{1, 2\}$ from above.

3 Application of Stephan's method to geodetic MBVP

The method of Stephan may be applied to any geodetic MBVP in order to transform it into a system of integral equations. Then, we have to study the solvability of the system, making use of the procedure as sketched above, and have to investigate the equivalence of the geodetic MBVP with the system of integral equations. For instance, let us consider the Dirichlet-oblique MBVP

$$\begin{aligned} \Delta u &= 0 && \text{in } D^c \\ u &= g_o && \text{on } S_o \\ D_\tau u &= g_c && \text{on } S_c \\ u &= O(|x|^{-1}), && |x| \rightarrow \infty, \end{aligned} \quad (22)$$

with Dirichlet data on the oceans and oblique-derivative data on the continents. γ_o and γ_c denotes the restriction to S_o and S_c , respectively. τ defines an oblique unit vector field on S , pointing into D . This problem has been studied by Keller (1996). It results after linearization of the non-linear fixed altimetry-gravimetry MBVP, which assumes that the geometry of

the Earth's surface is known and that gravity potential and gravity is given in ocean areas and continental areas, respectively. Keller (1996) shows the existence and uniqueness of the solution using the Kelvin transformation and the Lax-Milgram theorem. In order to transform the Dirichlet-oblique MBVP into an integral equation, we first need a representation formula that connects the Cauchy-data $\alpha := \gamma u$ and $\beta := \frac{1}{\langle n, \tau \rangle} \gamma(D_\tau u)$ with the unknown function u . Starting from Green's third identity (6), we obtain

$$u(x) = \int_S \frac{\alpha(y)}{\langle n, \tau \rangle(y)} D_r s(x-y) dS(y) - \int_S s(x-y) \beta(y) dS(y) + \int_S \langle \epsilon, \nabla(\alpha s) \rangle(y) dS(y), \quad x \in D^c, \quad (23)$$

with the unit vectors $r = 2\langle n, \tau \rangle n - \tau$ and $\epsilon = \frac{\tau}{\langle n, \tau \rangle} - n$. Since $\nabla(\alpha s) = \text{Grad}(\alpha s) + n D_n(\alpha s)$, and observing that $\langle \epsilon, n \rangle = 0$, we obtain $\langle \epsilon, \nabla(\alpha s) \rangle = \langle \epsilon, \text{Grad}(\alpha s) \rangle$. Grad denotes the surface gradient operator. Moreover, since S is a closed surface, it holds

$$\int_S \langle \epsilon, \text{Grad}(\alpha s) \rangle dS = - \int_S \alpha s \text{Div} \epsilon dS, \quad (24)$$

with the surface divergence operator Div . Therefore, we obtain for (23)

$$u(x) = \int_S \frac{\alpha(y)}{\langle n, \tau \rangle(y)} D_r s(x-y) dS(y) - \int_S s(x-y) \beta(y) dS(y) - \int_S \alpha(y) s(x-y) (\text{Div} \epsilon)(y) dS(y), \quad x \in D^c. \quad (25)$$

Equation (25) is our new representation formula. It is called "generalized Green-identity" (cf. Klees (1992, 1997)). Defining the oblique-derivative differential operator

$$P_r := \frac{1}{\langle n, \tau \rangle} D_r - \text{Div} \epsilon I, \quad (26)$$

we obtain the final form of our representation formula:

$$u(x) = \int_S P_{r(y)} s(x-y) \alpha(y) dS(y) - \int_S s(x-y) \beta(y) dS(y), \quad x \in D^c. \quad (27)$$

Observing the jump relations for the single layer potential and its gradient, we obtain the boundary integral equation (cf. Klees (1997))

$$\frac{1}{2} u(x) = \int_S P_{r(y)} s(x-y) \alpha(y) dS(y) - \int_S \beta(y) s(x-y) dS(y), \quad x \in S. \quad (28)$$

Taking the oblique derivative of (27), we obtain for the limit to the boundary

$$\frac{1}{2} \delta_1 \beta + \frac{1}{2} \delta_2 \alpha = \int_S D_{\tau(x)} P_{r(y)} s(x-y) \alpha(y) dS(y) - \int_S \beta(y) D_{\tau(x)} s(x-y) dS(y), \quad x \in S, \quad (29)$$

with

$$\delta_1 := \langle n, \tau \rangle - 1$$

and

$$\delta_2 := - \left(\langle n, \tau \rangle \text{Div} \epsilon - \frac{D_\tau \langle n, \tau \rangle}{|\langle n, \tau \rangle|^2} \right).$$

Defining the new boundary integral operators

$$\begin{aligned} (X\chi)(x) &:= \int_S P_{\tau(y)} s(x-y) \chi(y) dS(y), \\ x &\in S, \\ (U\chi)(x) &:= \int_S D_{\tau(x)} P_{\tau(y)} s(x-y) \chi(y) dS(y), \\ x &\in S, \\ (W\chi)(x) &:= \int_S D_{\tau(x)} s(x-y) \chi(y) dS(y), \\ x &\in S, \end{aligned} \quad (30)$$

we obtain the following system of integral equations by restricting (28) first to S_o and then to S_c :

$$\begin{pmatrix} X_{co} & -V_{oo} \\ \frac{1}{2}I - X_{cc} & V_{oc} \end{pmatrix} \begin{pmatrix} \alpha \\ \beta \end{pmatrix} = \begin{pmatrix} \frac{1}{2}I - X_{oo} & V_{co} \\ X_{oc} & -V_{cc} \end{pmatrix} \begin{pmatrix} g_o \\ \frac{1}{\langle n, \tau \rangle} g_c \end{pmatrix}. \quad (31)$$

Equation (31) defines a mixed system of boundary integral equations; it is of the second kind w.r.t. α and of the first kind w.r.t. β . We can derive alternative integral equations, e.g., by restricting (28) to S_o and (29) to S_c and vice versa. For instance, restricting (28) to S_c and (29) to S_o , we obtain a system of second kind integral equations for the unknowns α and β :

$$\begin{pmatrix} U_{co} & -W_{oo} - \frac{1}{2}\langle n, \tau \rangle \delta_1 I \\ \frac{1}{2}I - X_{cc} & V_{oc} \end{pmatrix} \begin{pmatrix} \alpha \\ \beta \end{pmatrix} = \begin{pmatrix} \frac{1}{2}\delta_2 I - U_{oo} & W_{co} \\ X_{oc} & -V_{cc} \end{pmatrix} \begin{pmatrix} g_o \\ \frac{1}{\langle n, \tau \rangle} g_c \end{pmatrix}, \quad (32)$$

Analogously, restricting (28) to S_o and (29) to S_c , we obtain the mixed system of integral equations:

$$\begin{pmatrix} \frac{1}{2}\delta_2 I - U_{cc} & W_{oc} \\ X_{co} & -V_{oo} \end{pmatrix} \begin{pmatrix} \alpha \\ \beta \end{pmatrix} = \begin{pmatrix} U_{oc} & -\frac{1}{2}\langle n, \tau \rangle \delta_1 I - W_{cc} \\ \frac{1}{2}I - X_{oo} & V_{co} \end{pmatrix} \begin{pmatrix} g_o \\ \frac{1}{\langle n, \tau \rangle} g_c \end{pmatrix}. \quad (33)$$

Like (31) it is of the second kind w.r.t. α and of the first kind w.r.t. β . The boundary operators in (31), (32), and (33) have the following mapping properties: For some real s , depending on the smoothness of S , and $i, k = \{o, c\}$, the mappings

$$\begin{aligned} X_{ik} &: \tilde{H}^s(S_i) \rightarrow H^s(S_k), \\ U_{ik} &: \tilde{H}^{s+1}(S_i) \rightarrow H^s(S_k), \\ W_{ik} &: \tilde{H}^s(S_i) \rightarrow H^s(S_k) \end{aligned} \quad (34)$$

are continuous. X_{ik} and W_{ik} define strongly singular integral operators, which are pseudodifferential operators of order 0; U_{ik} is a hypersingular integral operator, which has order 1. For the property of V_{ik} , see (9). What still has to be done is to investigate the solvability of the systems (31)-(33) and to prove the equivalence of the Dirichlet-oblique MBVP with the systems of integral equations. This can be done following the procedure of Stephan (1987).

4 The method of Sansò and Stock

Sansò and Stock (1985) consider the Robin-Neumann mixed boundary value problem

$$\begin{aligned} \Delta u &= 0 && \text{in } D^c \\ D_n u &= g_o && \text{on } S_o \\ D_n u + \frac{2}{R}u &= g_c && \text{on } S_c \\ u &= O(|x|^{-1}), && |x| \rightarrow \infty, \end{aligned} \quad (35)$$

where S is the surface of a sphere with radius R . They look for a solution $u \in H_{loc}^\lambda(S)$ for given boundary data $g_o \in H^{\lambda-1}(S_o)$ and $g_c \in H^{\lambda-1}(S_c)$ with $\frac{1}{2} < \lambda < \frac{3}{2}$. The transformation into a boundary integral equation is based on the explicit solution of the Neumann BVP for a sphere, which is known as Hotine's formula:

$$u(x) = -\frac{R}{4\pi} \int_S H(x-y) (\gamma D_n u)(y) dS(y), \quad x \in S, \quad (36)$$

with the Green function of the second kind (Hotine function, Neumann function)

$$H(x-y) = \frac{2R}{|x-y|} - \ln \left(1 + \frac{2R}{|x-y|} \right). \quad (37)$$

Defining the integral operator

$$E\chi(x) := -\frac{R}{4\pi} \int_S H(x-y) \chi(y) dS(y), \quad x \in S, \quad (38)$$

equation (36) can be written as $\mu = E\nu$. As in Section 2, $\left(\begin{smallmatrix} \mu \\ \nu \end{smallmatrix}\right)$ define the traces on S $\left(\begin{smallmatrix} \gamma u \\ \gamma D_n u \end{smallmatrix}\right)$. Observing the boundary condition of the Robin-Neumann MBVP, we have

$$\mu = E\nu = E_{oS}\nu + E_{cS}\nu = E_{oS}g_o + E_{cS} \left(g_c - \frac{2}{R}\mu \right), \quad (39)$$

hence,

$$\left(I + \frac{2}{R}E_{cS} \right) \mu = E_{oS}g_o. \quad (40)$$

Equation (40) is an integral equation of the second kind for the unknown Cauchy-data μ . Substituting $\mu = \mu_0 + \frac{R}{2}lg_o$, we obtain

$$A\mu_0 := \left(I + \frac{2}{R}E_{cS} \right) \mu_0 = - \left(\frac{R}{2}I - E \right) lg_o =: Blg_o. \quad (41)$$

The operators $A, B : H^s(S) \rightarrow H^s(S)$ are continuous for some real s depending on the smoothness of S . Analogously to Section 2, we can prove the solvability of (41) and the equivalence of the integral equation with the Robin-Neumann MBVP. We omit the details. In order to solve the integral equation, we can apply, e.g., the Nyström method, collocation boundary element methods or Galerkin boundary element methods.

The method relies on the spherical topology since only then the Neumann function is known, i.e., there is an explicit solution of the Neumann problem available. However, the basic idea can easily be generalized if instead of Hotine's formula Green's third identity is used. With the results of Section 2, we have

$$\frac{1}{2}\mu = K\mu - V\nu = K\mu - V_{oS}\nu - V_{cS}\nu = K\mu - V_{oS}\nu - V_{cS}\left(g_c - \frac{2}{R}\mu\right), \quad (42)$$

hence,

$$\left(\frac{1}{2}I - K - \frac{2}{R}V_{cS}\right)\mu = -V_{oS}g_o - V_{cS}g_c. \quad (43)$$

The operator $A := \left(\frac{1}{2}I - K - \frac{2}{R}V_{cS}\right)$ is continuous from $H^s(S) \rightarrow H^s(S)$, the operator $B := -\left(\frac{V_{oS}}{V_{cS}}\right)^T$ is continuous from $H^{s-1}(S_o) \times H^{s-1}(S_c) \rightarrow H^s(S)$.

Alternative integral equations can be derived in different ways making use of (8) and restricting to S_o or S_c . If on S_c $\gamma_c\nu$ is replaced by $g_c - \frac{2}{R}\gamma_c\mu$, we obtain only one integral equation for the unknown $\mu = \gamma u$. For instance, when restricting the first equation in (8) to S_o , we obtain the second kind integral equation

$$\left(\frac{1}{2}I - K_{S_o} - \frac{2}{R}V_{c_o}\right)\mu = -V_{o_o}g_o - V_{c_o}g_c, \quad x \in S_o, \quad (44)$$

with weakly singular kernels. When restricting the second equation in (8) to S_c , we obtain a second kind integral equation with weakly singular and hypersingular kernels

$$\left(\frac{1}{2}I + \frac{2}{R}K'_{cc} - D_{S_c}\right)\mu = \left(\frac{1}{2}I + K'_{cc}\right)g_c + K'_{oc}g_o, \quad x \in S_c. \quad (45)$$

We can also derive a first order integral equation by restricting the second equation in (8) to S_o :

$$\left(D_{S_o} - \frac{2}{R}K'_{c_o}\right)\mu = -\left(\frac{1}{2}I + K'_{o_o}\right)g_o + K'_{c_o}g_c, \quad x \in S_o. \quad (46)$$

The corresponding kernels are weakly singular and hypersingular. What remains is to prove the solvability of the integral equations and the equivalence with the original MBVP. The prove can easily be done using the procedure and results of Section 2. We omit the details.

5 The method of Mayer

Mayer (1997) considers the altimetry-gravimetry II MBVP in non-linear form. There are two sources of nonlinearities:

1. gravity is a non-linear functional of the potential, and
2. the boundary surface is partly free (over the continents, unobservable by altimeter radar).

The new idea of Mayer is to perform a linearization with respect to source 1 (gravity) only, and later on, to solve the resulting (still non-linear) problem by a special iteration procedure. This approach is justified by recent findings of Heck and Seitz (1993), that source 2 (free boundary) is the more severe source of non-linearity in geodetic boundary value problems. Consequently, if any iteration will be necessary, then w.r.t. source 2.

The formulation of the partly linearized altimetry-gravimetry II MBVP is:

$$\begin{aligned}
 \Delta u &= 0 && \text{in } D^c \\
 D_\tau u &= g && \text{on } S \\
 u &= g_c && \text{on } S_c \\
 u &= O(|x|^{-1}), && |x| \rightarrow \infty.
 \end{aligned} \tag{47}$$

The oceanic surface $S_o := S \setminus S_c$ is assumed to be known, as well. If in addition S_c were known instead of g_c , the resulting BVP would be identical to the linearized fixed gravimetric BVP, which in turn equals the classical oblique BVP for the Laplace equation. This problem is much easier to solve because it is *not mixed*. Some theoretical results exist (e.g., Klees (1992)), and have been augmented recently by Mayer (1997). Also from the numerical point of view, this problem is solvable using boundary element methods (e.g., Klees (1992)).

A boundary integral equation for the linearized fixed gravimetric BVP is derived from a representation formula, for which Mayer prefers a combined double- and single-layer potential

$$u(x) = (K\chi)(x) + \kappa(V\chi)(x), \quad x \in D^c, \tag{48}$$

where κ is an arbitrary positive real number, and χ is the surface density. Defining the operators

$$\begin{aligned}
 (Y\chi)(x) &:= \int_S D_{\tau(x)} D_{n(y)} s(x-y) \chi(y) dS(y), \\
 &x \in S \\
 (Z\chi)(x) &:= \frac{1}{2} \langle \text{Grad}\chi, \tau \rangle,
 \end{aligned} \tag{49}$$

we obtain an integro-differential equation of the second kind on S

$$A\chi := \left(-\frac{1}{2} \kappa \langle n, \tau \rangle I + Z + Y + \kappa W \right) \chi = g. \tag{50}$$

Formally, the solution of this integro-differential equation can be written as

$$\chi = A^{-1}g. \tag{51}$$

Hence, the desired potential function is

$$u(x) = ([K + \kappa V]A^{-1}g)(x), \quad x \in D^c. \tag{52}$$

Now, we return to the actual problem (47), where additionally the Dirichlet condition over the continents

$$\gamma_c u = g_c \quad (53)$$

has to be fulfilled. The only unknown of this equation is the boundary surface S_c . Therefore, we end up with an operator equation, which must be solved for S_c :

$$Q(S_c, g) := \gamma_c u = \gamma_c \{ (K + \kappa V) A^{-1} g \} = g_c, \quad \text{on } S_c \quad (54)$$

The operator $Q(S_c, g)$ is *non-linear* in the first argument. The complicated structure of this operator is the price we have to pay for the striking simplicity of the first solution step (51),(52). The nonlinear operator equation (54) must be solved iteratively, starting from an initial guess \tilde{S}_c for S_c , which in geodesy is known as the telluroid. However, note that $\tilde{S} := \tilde{S}_c \cup S_o$ must be a *closed* surface, which is not guaranteed by classical definitions of the telluroid. Quite formally, let us express the *Fréchet expansion* of Q as

$$Q(S_c, g) = Q(\tilde{S}_c, g) + \frac{\partial Q}{\partial S_c}(\tilde{S}_c, g)(S_c - \tilde{S}_c) + o(|S_c - \tilde{S}_c|). \quad (55)$$

This expansion suggests an iterative procedure of Newton type: Let $S_c^{(0)} := \tilde{S}_c$; for $n = 1, 2, \dots$:

$$S_c^{(n)} := S_c^{(n-1)} + \left(\frac{\partial Q}{\partial S_c}(S_c^{(n-1)}, g) \right)^{-1} (g_c - Q(S_c^{(n-1)}, g)). \quad (56)$$

So far, nothing can be said about the feasibility of this suggestion, neither about the *existence and uniqueness* of the inverse Fréchet derivative nor about the *convergence* of this procedure. Moreover, due to the complicated structure of Q there is even less hope to obtain similar results as with classical geodetic approaches. Mayer (1997) has also derived an *explicit expression* for the Fréchet derivative of Q . However, the complexity of this expression will certainly prevent any practical application in geodesy.

Nonetheless, Mayer (1997) has shown that there always exist interesting alternatives to the standard geodetic techniques.

References

- Arnold, K. (1981). Complex evaluation of gravity anomalies and data obtained from satellite altimetry. *Gerlands Beitr. Geophys.*, **90**, 38–42.
- Arnold, K. (1984). The compatibility conditions, the uniqueness and the solution of the mixed boundary value problem in geodesy. *Gerlands Beitr. Geophys.*, **93**, 339–355.
- Bjerhammar, A. (1983). A stochastic approach to the mixed boundary value problem in physical geodesy. *Geodesy in Transition (dedicated to H. Moritz)*, pages 25–42.
- Heck, B. and Seitz, K. (1993). Effects of non-linearity in the geodetic boundary value problems. , Veröffentl. d. Dt. Geod. Kommiss. b. d. Bayer. Akad. d. Wiss., Reihe A, Heft Nr. 109, München.
- Hofmann-Wellenhof, B. (1985). The use of multipoles for the altimetry-gravimetry problem. Internal report, Technical University Graz.

- Holota, P. (1982). The altimetry-gravimetry boundary value problem. In *Utilization of observations of artificial satellites of the earth for the purposes of geodesy*, Proc. Int. Sci. Conf. of Sec. 6, Intercosmos, Bulgaria, Albena, Varna, Sept. 1980, pages 243–249.
- Holota, P. (1983a). The altimetry gravimetry boundary value problem i: linearization, friedrich's inequality. *Bolletino de Geodesia e Scienze Affini*, **17**, 13–32.
- Holota, P. (1983b). The altimetry gravimetry boundary value problem ii: weak solution, v-ellipticity. *Bolletino de Geodesia e Scienze Affini*, **17**, 69–84.
- Keller, W. (1996). On the scalar fixed altimetry gravimetry boundary value problem. *Journal of Geodesy*, **70**, 459–469.
- Klees, R. (1992). *Lösung des fixen geodätischen Randwertproblems mit Hilfe der Randelementmethode*. Ph.D. thesis, Deutsche Geodätische Kommission, Reihe C, No. 382, München.
- Klees, R. (1997). Topics of boundary element methods. In F. Sanso and R. Rummel, editors, *Geodetic Boundary Value Problems in View of the One Centimeter Geoid*, volume 65 of *Lecture Notes in Earth Sciences*, pages 482–531. Springer.
- Mainville, A. (1986). The altimetry-gravimetry problem using orthonormal base functions. 373, Department of Geodetic Science and Surveying, The Ohio State University.
- Mayer, J. (1997). *Zur Lösung von geodätischen Randwertproblemen durch einen hypersingulären Potentialansatz*. Ph.D. thesis, Universität Kiel.
- Miranda, C. (1970). *Partial Differential Equations of Elliptic Type*. Springer.
- Sacerdote, F. and Sansò, F. (1983a). A contribution to the analysis of altimetry-gravimetry problems. *Bulletin Geodesique*, **57**, 257–272.
- Sacerdote, F. and Sansò, F. (1983b). A contribution to the analysis of the altimetry-gravimetry problem. In *Figure of the Earth, the Moon and other Planets*, Monography Series of VUGTK, pages 123–139.
- Sacerdote, F. and Sansò, F. (1987). Further remarks on the altimetry gravimetry problems. *Bulletin Geodesique*, **61**, 65–82.
- Sansò, F. (1993). Theory of geodetic b.v.p.s applied to the analysis of altimetric data. In R. Rummel and F. Sanso, editors, *Satellite Altimetry in Geodesy and Oceanography*, volume 50 of *Lecture Notes in Earth Sciences*, pages 318–373. Springer.
- Sansò, F. and Stock, B. (1985). A numerical experiment in the altimetry-gravimetry problem ii. *Manuscripta Geodaetica*, **10**, 23–31.
- Sjöberg, L. (1982). On the altimetry-gravimetry boundary problem. *Bollettino di Geodesia e Scienze Affini*, **4**, 377–392.
- Stephan, E. (1987). Boundary integral equations for mixed boundary value problems in \mathbb{R}^3 . *Math. Nachr.*, **134**, 21–53.
- Svensson, L. (1983). Solution of the altimetry-gravimetry problem. *Bulletin Geodesique*, **57**, 332–353.
- Svensson, L. (1988). Some remarks on the altimetry-gravimetry problem. *Manuscripta Geodaetica*, **13**, 63–74.
- von Petersdorff, T. and Stephan, E. (1990). Regularity of mixed boundary value problems in \mathbb{R}^3 and boundary element methods on graded meshes. *Mathematical Methods in the Applied Sciences*, **12**, 229–249.

1. The first part of the document is a list of names and addresses of the members of the committee.

2. The second part is a list of the names and addresses of the members of the committee.

3. The third part is a list of the names and addresses of the members of the committee.

4. The fourth part is a list of the names and addresses of the members of the committee.

5. The fifth part is a list of the names and addresses of the members of the committee.

6. The sixth part is a list of the names and addresses of the members of the committee.

7. The seventh part is a list of the names and addresses of the members of the committee.

8. The eighth part is a list of the names and addresses of the members of the committee.

9. The ninth part is a list of the names and addresses of the members of the committee.

10. The tenth part is a list of the names and addresses of the members of the committee.

11. The eleventh part is a list of the names and addresses of the members of the committee.

12. The twelfth part is a list of the names and addresses of the members of the committee.

13. The thirteenth part is a list of the names and addresses of the members of the committee.

14. The fourteenth part is a list of the names and addresses of the members of the committee.

15. The fifteenth part is a list of the names and addresses of the members of the committee.

16. The sixteenth part is a list of the names and addresses of the members of the committee.

17. The seventeenth part is a list of the names and addresses of the members of the committee.

18. The eighteenth part is a list of the names and addresses of the members of the committee.

On Singular Surface Integrals in Physical Geodesy

Roland Klees

Abstract

We consider integrals over closed surfaces with non-integrable point singularities that arise in Physical Geodesy. We investigate their behavior under smooth parameter transformations and show how they can be regularized making use of the definition of singular integrals. Finally, we address the problem of numerical integration of the regularized integrals.

1 Introduction

The objective of the paper is to study integrals of type

$$\int_U \tilde{K}(x, y - x) \tilde{\mu}(y) dU(y), \quad x \in U. \quad (1)$$

U is a piecewise smooth closed surface, and y and x are points on U ; $K(x, y - x)$ is the kernel function, which has order $O(|y - x|^{-s})$ as $y - x \rightarrow 0$. $\mu(y)$ is the density function, which is assumed to be sufficiently smooth, and s is the order of the singularity. If $s = 1$ the kernel K is "weakly singular", but still absolutely integrable. Examples are the single layer potential, the double layer potential, and the formulas of Hotine and Stokes. Cubature formulas for weakly singular surface integrals have been studied by, e.g. Klees (1996). If $s = 2$ and $s = 3$ the kernel K is called "strongly singular" and "hypersingular", respectively, and the integral is not absolutely integrable. For instance the oblique derivative of the single layer potential (τ is the oblique unit vector)

$$\int_U \frac{\langle \tau(x), y - x \rangle}{|y - x|^3} \tilde{\mu}(y) dU(y), \quad (2)$$

the integral of Vening-Meinesz (e is an oblique unit vector, St the Stokes function)

$$\int_U \langle e(x), \nabla_x St(y - x) \rangle \tilde{\mu}(y) dU(y), \quad (3)$$

and Molodensky's G_1 -Term

$$\int_U \frac{h(x) - h(y)}{|y - x|^3} \tilde{\mu}(y) dU(y) \quad (4)$$

have kernel functions of order $O(|y-x|^{-2})$, i.e. $s = 2$. The oblique derivative of the double layer potential (n is the normal unit vector on U)

$$\int_U \left(\frac{\langle n(y), \tau(x) \rangle}{|y-x|^3} - 3 \frac{\langle n(y), y-x \rangle \langle \tau(x), y-x \rangle}{|y-x|^5} \right) \tilde{\mu}(y) dU(y) \quad (5)$$

has a kernel function with order of singularity $s = 3$, and the multipole of order p

$$\int_U \frac{\partial^p}{\partial n(y)^p} \left[\frac{1}{|y-x|} \right] \tilde{\mu}(y) dU(y) \quad (6)$$

has order $s = p$. The paper aims at the efficient numerical computation of integrals with non-integrable kernels, i.e. when $s \geq 2$. Section 2 is devoted to the behavior under smooth parameter transformations of singular surface integrals. In Section 3 we will regularize the integrals making use of the definition of singular surface integrals and of polar coordinates in the parameter domain. In Section 4, we apply the theory to a numerical example and discuss some aspects of the numerical computation of the regularized integrals. Section 5, finally, contains a short summary and some conclusions.

The motivation for this paper is three-fold. Firstly, integrals with point singularities of potential type that are not absolutely integrable have not been discussed yet systematically in geodetic literature. There are only a few references which cover some very special cases of (1), mostly for planar boundary surfaces U , e.g. Bosch (1977), Shaofeng and Xurong (1991), Bian and Sum (1994). Secondly, a change of variables (parameter transformation) is commonly performed in singular integrals without taking care of the additional terms that may appear since in general the substitution rule as known for absolutely integrable kernels does not apply to singular kernels (e.g. Vijayakumar and Cormack (1988)). Thirdly, singular integrals are avoided as often as possible in Physical Geodesy, because it is believed that they cannot be computed efficiently. A lot of effort is put into the regularization making use of standard theorems of vector analysis on surfaces. This results in absolutely integrable kernels which often have very complicated structures making the numerical computation more elaborate than the strategy to be discussed in this paper (e.g. Meissl (1971), Hofmann-Wellenhof (1983)).

A proper analysis of strongly singular and hypersingular integrals has been done by Kieser (1991), and Schwab and Wendland (1992a). Numerical integration formulas, which properly take into account the behavior of singular integrals under smooth parameter transformations, have been developed in, e.g., Guiggiani (1991), Klees (1992), Schwab and Wendland (1992b), Kieser *et al.* (1992), Guiggiani *et al.* (1992).

2 Parameter transformation

Let us first study the behavior of singular integrals under smooth parameter transformations. For any fixed computation point $x \in U$ we may split-up the integral over U into two parts, the first taken over $U \setminus S$ and the second taken over S , where S denotes some neighborhood of x . The integral over $U \setminus S$ is regular since $x \notin U \setminus S$, so $|y-x|$ cannot become zero. Therefore, we can limit to the integral over S . For convenience we assume S to be the image of the standard

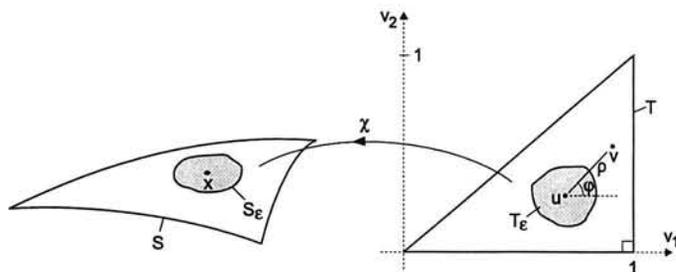


Fig. 1. Surface patch, reference triangle, and mapping χ

triangle

$$T = \{v : 0 \leq v_1 \leq 1, 0 \leq v_2 \leq v_1\} \quad (7)$$

under a smooth mapping χ , i.e. $S = \chi(T)$ (cf. Figure 1). Therefore, we consider the integral

$$\int_S \tilde{K}(x, y) \tilde{\mu}(y) dS(y), \quad x \in S, \quad (8)$$

where S is a triangular surface patch containing the computation point x in its interior. In order to transform the integral (8) into an integral over T let us first exclude a neighborhood of the computation point, say, S_ε . S_ε may be any domain containing x that shrinks down to x as $\varepsilon \rightarrow 0$. Moreover we assume S_ε to be star-shaped w.r.t. x (cf. Figure 1). Then, instead of (8), we consider the *regular* integral

$$\int_{S \setminus S_\varepsilon} \tilde{K}(x, y) \tilde{\mu}(y) dS(y), \quad x \in S_\varepsilon. \quad (9)$$

Later, we will investigate the behavior of this integral as $\varepsilon \rightarrow 0$, but for the moment we assume $\varepsilon > \varepsilon_0$ with some positive number ε_0 . Now, since the integral is regular, we can perform the parameter transformation applying the usual substitution rule to the regular integral (9): With $x = \chi(u)$, $y = \chi(v)$, $\tilde{K}(x, y - x) = K(u, v - u)$, and $\tilde{\mu}(y) dS(y) = \mu(v) dT(v)$, we obtain

$$\int_{S \setminus S_\varepsilon} \tilde{K}(x, y - x) \tilde{\mu}(y) dS(y) = \int_{T \setminus T_\varepsilon} K(u, v - u) \mu(v) dT(v), \quad (10)$$

where $T_\varepsilon = \chi^{-1}(S_\varepsilon)$. The kernel $K(u, v - u)$ and the density $\mu(v)$ can be expanded into a Taylor series around u at $\rho = 0$, where ρ, φ are the polar coordinates w.r.t. u . We obtain

$$K(u, v - u) = \sum_{k=0}^{s-2} \rho^{k-s} K_k(u, \varphi) + (R_{s-2}K)(u, v), \quad (11)$$

and

$$\mu(v) = \sum_{j=0}^{s-2} \rho^j \mu_j(\varphi) + (R_{s-2}\mu)(v). \quad (12)$$

The remainders $(R_{s-2}K)(u, v)$ and $(R_{s-2}\mu)(v)$ are of the order $O(\rho^{-1})$ and $O(\rho^{s-1})$, respectively, as $\rho \rightarrow 0$. Thus,

$$\begin{aligned} & \int_{T \setminus T_\varepsilon} K(u, v - u) \mu(v) dT(v) \\ &= \int_{T \setminus T_\varepsilon} K(u, v - u) (R_{s-2}\mu)(v) dT(v) + \sum_{j=0}^{s-2} \int_{T \setminus T_\varepsilon} (R_{s-2}K)(u, v) \rho^j \mu_j(\varphi) dT(v) \\ & \quad + \sum_{-1 \leq i+j-s \leq s-4} \int_{T \setminus T_\varepsilon} \rho^{i+j-s} K_i(u, \varphi) \mu_j(\varphi) dT(v) \\ & \quad + \underbrace{\sum_{i=0}^{s-2} \int_{T \setminus T_\varepsilon} \rho^{i-s} g_i(u, \varphi) dT}_{s-1 \text{ terms with kernels } O(\rho^{-\sigma}), 2 \leq \sigma \leq s}, \end{aligned} \quad (13)$$

with

$$g_i(u, \varphi) := \sum_{j=0}^i K_{i-j}(u, \varphi) \mu_j(\varphi). \quad (14)$$

The first three integrals on the right-hand side of (13) converge absolutely to the integral over T as $\varepsilon \rightarrow 0$. They have kernels $O(\rho^\sigma)$, $\sigma \geq -1$, i.e. the limit $\varepsilon \rightarrow 0$ exists, as regular or improper integral. Therefore, we can limit to the last term on the right-hand side of (13), which consists of terms $O(\rho^{-\sigma})$, $2 \leq \sigma \leq s$. In order to simplify the notation, we will omit from now on the dependency of the functions of the computation point u , i.e. we write e.g. instead of $g(u, \varphi)$ simply $g(\varphi)$ etc. The terms with kernels of order $O(\rho^{-\sigma})$, $2 \leq \sigma \leq s$, i.e.

$$F_i(\varepsilon) := \int_{T \setminus T_\varepsilon} \rho^{i-s} g_i(\varphi) dT, \quad i = 0, \dots, s-2, \quad (15)$$

can be integrated analytically w.r.t. the variable ρ . Assuming that $R(\varphi)$ is the boundary of the standard triangle T in polar coordinates, and $R_\varepsilon(\varphi)$ is the boundary of T_ε in polar coordinates (cf. Figure 2), we obtain, observing $dT = \rho d\rho d\varphi$,

$$\begin{aligned} F_i(\varepsilon) &= \int_0^{2\pi} g_i(\varphi) \int_{R_\varepsilon(\varphi)}^{R(\varphi)} \rho^{i-s+1} d\rho d\varphi \\ &= \int_0^{2\pi} g_i(\varphi) \left\{ \begin{array}{ll} \ln R(\varphi) - \ln R_\varepsilon(\varphi), & i = s-2 \\ \frac{1}{m} [R_\varepsilon^{-m}(\varphi) - R^{-m}(\varphi)], & 0 \leq i < s-2 \end{array} \right\} d\varphi, \end{aligned} \quad (16)$$

where we have introduced the new variable $m = s-2-i$ if $i \neq s-2$, i.e. for m it holds $0 < m \leq s-2$.

3 Regularization

In order to study what happens as $\varepsilon \rightarrow 0$, we need to know the equation of the boundary of T_ε , $\rho = R_\varepsilon(u, \varphi)$, in terms of ε . Since we assumed that T_ε is star-shaped w.r.t. u and that T_ε shrinks

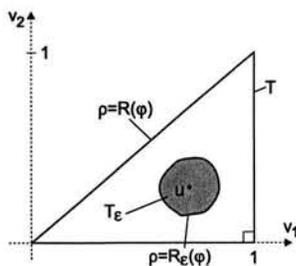


Fig. 2. ε -neighborhood T_ε and boundary of T in local coordinates

down to zero as $\varepsilon \rightarrow 0$, $R_\varepsilon(\varphi)$ has in general an expansion of the form

$$R_\varepsilon(\varphi) = \varepsilon \sum_{i \geq 0} d_i(\varphi) \varepsilon^i, \quad \varepsilon \rightarrow 0. \quad (17)$$

Therefore,

$$\ln R_\varepsilon(\varphi) = \ln \varepsilon + \ln d_0(\varphi) + o(1), \quad \varepsilon \rightarrow 0, \quad (18)$$

and

$$R_\varepsilon^{-m}(\varphi) = d_{m,m}(\varphi) + \sum_{n=0}^{m-1} d_{n,m}(\varphi) \varepsilon^{n-m} + o(1), \quad \varepsilon \rightarrow 0. \quad (19)$$

The functions $d_{n,m}(\varphi)$ can be expressed in terms of the functions $d_i(\varphi)$. Inserting (18), (19) into (16) yields

$$F_i(\varepsilon) \sim \int_0^{2\pi} g_i(\varphi) [\ln R(\varphi) - \ln d_0(\varphi)] d\varphi - \underbrace{\ln \varepsilon \int_0^{2\pi} g_i(\varphi) d\varphi}_{\text{unbounded term, } \varepsilon \rightarrow 0}, \quad i = s-2, \quad (20)$$

and

$$F_i(\varepsilon) \sim \frac{1}{m} \int_0^{2\pi} g_i(\varphi) [d_{m,m}(\varphi) - R^{-m}(\varphi)] d\varphi + \underbrace{\frac{1}{m} \sum_{n=0}^{m-1} \varepsilon^{n-m} \int_0^{2\pi} g_i(\varphi) d_{n,m}(\varphi) d\varphi}_{\text{unbounded terms } O(\varepsilon^{-1}), \dots, O(\varepsilon^{-m}), \varepsilon \rightarrow 0} \\ 0 \leq i < s-2. \quad (21)$$

In general, the terms with $\ln \varepsilon, \varepsilon^{-1}, \dots, \varepsilon^{-m}$ are unbounded in the limit $\varepsilon \rightarrow 0$, i.e. $\lim_{\varepsilon \rightarrow 0} F_i(\varepsilon)$ does not exist. Neglecting all unbounded terms yields the so-called *finite part* of $F_i(\varepsilon)$, written f.p. $F_i(\varepsilon)$:

$$\text{f.p. } F_i(\varepsilon) = \int_0^{2\pi} g_i(\varphi) [\ln R(\varphi) - \ln d_0(\varphi)] d\varphi, \quad i = s-2, \quad (22)$$

and

$$\text{f.p. } F_i(\varepsilon) = \frac{1}{m} \int_0^{2\pi} g_i(\varphi) [d_{m,m}(\varphi) - R^{-m}(\varphi)] d\varphi, \quad 0 \leq i < s-2. \quad (23)$$

The limit $\lim_{\varepsilon \rightarrow 0} F_i(\varepsilon)$ exists if and only if

$$\int_0^{2\pi} g_i(\varphi) d\varphi = 0, \quad i = s-2, \quad (24)$$

and

$$\int_0^{2\pi} g_i(\varphi) d_{n,m}(\varphi) d\varphi = 0, \quad \text{for all } n = 0, \dots, m-1, \quad 0 \leq i < s-2, \quad (25)$$

respectively. Then it is called *Cauchy principal value* of $F_i(\varepsilon)$, written p.v. $F_i(\varepsilon)$, i.e.

$$\text{p.v. } F_i(\varepsilon) = \lim_{\varepsilon \rightarrow 0} F_i(\varepsilon) \quad (26)$$

if the limit exists. Note that the Cauchy principal value, if it exists, is equal to the finite part of the singular integral. The only difference is that the finite part of a singular integral imposes stronger requirements on the smoothness of the density μ .

The result has some implications for the choice of the neighborhood T_ε of the computation point u . Assume that there are two neighborhoods, $T_{1,\varepsilon}$ and $T_{2,\varepsilon}$, with star-shaped continuous boundaries $\rho = R_{1,\varepsilon}(\varphi)$ and $\rho = R_{2,\varepsilon}(\varphi)$, both shrinking down to zero as $\varepsilon \rightarrow 0$ (cf. Figure 3). Then, we may consider the difference

$$\Delta F_i(\varepsilon) = \int_{T_{2,\varepsilon} \setminus T_{1,\varepsilon}} \rho^{i-2} g_i(\varphi) dT, \quad (27)$$

and obtain for the finite part of $\Delta F_i(\varepsilon)$:

$$\text{f.p. } \Delta F_i(\varepsilon) = \begin{cases} \int_0^{2\pi} g_i(\varphi) [\ln d_0^{(2)}(\varphi) - \ln d_0^{(1)}(\varphi)] d\varphi, & i = s-2 \\ -\frac{1}{m} \int_0^{2\pi} g_i(\varphi) [d_{m,m}^{(2)}(\varphi) - d_{m,m}^{(1)}(\varphi)] d\varphi, & 0 \leq i < s-2 \end{cases} \quad (28)$$

Since f.p. $\Delta F_i(\varepsilon)$ is not equal to zero, we can conclude that the finite part depends on the shape of the ε -neighborhood of u . For instance, if $T_{1,\varepsilon}$ is the circle with radius ε and center u , and $T_{2,\varepsilon}$ is the ellipse with semi-axes ε and $\delta\varepsilon$, $\delta < 1$, centered at u , we obtain

$$R_{1,\varepsilon} = \varepsilon, \quad d_0^{(1)}(\varphi) = 1, \quad d_{m,m}^{(1)}(\varphi) = 0, \quad (29)$$

and

$$R_{2,\varepsilon} = \varepsilon \left[1 + \frac{1-\delta^2}{\delta^2} \sin^2 \varphi \right]^{-1/2},$$

$$d_0^{(2)}(\varphi) = \left[1 + \frac{1-\delta^2}{\delta^2} \sin^2 \varphi \right]^{-1/2}, \quad d_{m,m}^{(2)}(\varphi) = 0. \quad (30)$$

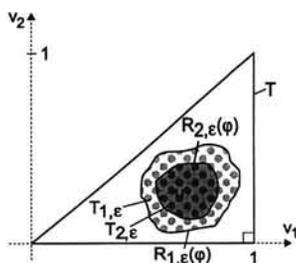


Fig. 3. Choice of the ε -neighborhood in local coordinates

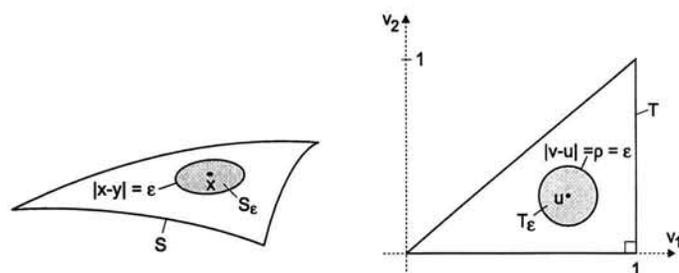


Fig. 4. "Identical" ε -neighborhoods in global and local coordinates

Thus,

$$\text{f.p. } \Delta F_i(\varepsilon) = \begin{cases} -\frac{1}{2} \int_0^{2\pi} g_i(\varphi) \ln \left[1 + \sin^2 \varphi \frac{1-\delta^2}{\delta^2} \right] d\varphi, & i = s-2 \\ 0, & 0 \leq i < s-2 \end{cases} \quad (31)$$

Therefore, for a given ε -neighborhood of the computation point x , S_ε , the choice of the neighborhood T_ε in the parameter domain is not arbitrary but given by $T_\varepsilon = \chi^{-1}(S_\varepsilon)$. In particular (cf. Figure 4)

$$\text{f.p. } \int_{S \setminus |y-x| < \varepsilon} \tilde{K}(x, y-x) \tilde{\mu}(y) dS(y) \neq \text{f.p. } \int_{T \setminus |v-u| < \varepsilon} K(u, v-u) \mu(v) dT(v). \quad (32)$$

Since it is always easier to work with "simple" neighborhoods, we may ask when the choice of the neighborhood T_ε is independent of the given neighborhood S_ε , i.e. when T_ε can be chosen arbitrarily, e.g. as the circle with radius ε . Kieser (1991) has proved that this only depends on the properties of the kernel function $\tilde{K}(x, y-x)$. Let $\tilde{K}(x, y-x)$ admit an expansion of the form

$$\tilde{K}(x, y-x) = \sum_{m=0}^{s-2} \mathcal{K}_m(x, z) + O(|z|^{-1}), \quad z := P_x(y-x), \quad (33)$$

where P_x is the orthogonal projection onto the tangent plane of S at x , (cf. Figure 5), and let the terms $\mathcal{K}_m(x, z)$ be homogeneous of degree $m-s$ w.r.t. the second variable, i.e.

$$\mathcal{K}_m(x, tz) = t^{m-s} \mathcal{K}_m(x, z), \quad t > 0. \quad (34)$$

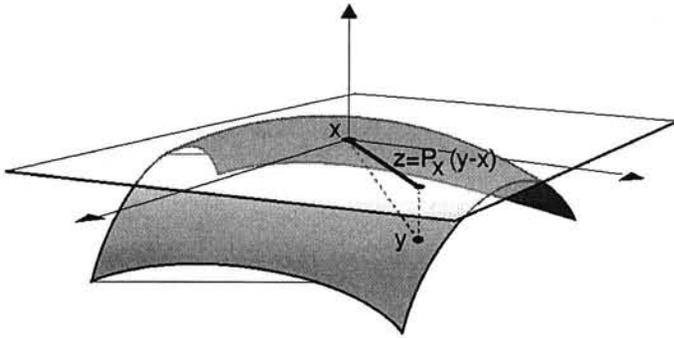


Fig. 5. Invariance of the finite part integral w.r.t. choice of the ε -neighborhood

If

$$\mathcal{K}_m(x, -z) = (-1)^{m+1-s} \mathcal{K}_m(x, z) \quad (35)$$

the finite part is independent of the ε -neighborhood T_ε , and the usual substitution rule can be applied. In particular it holds (cf. Figure 4)

$$\text{f.p.} \int_{S \setminus \{|y-x| < \varepsilon\}} \tilde{K}(x, y-x) \tilde{\mu}(y) dS(y) = \text{f.p.} \int_{T \setminus \{|v-u| < \varepsilon\}} K(u, v-u) \mu(v) dT(v). \quad (36)$$

For applications in Physical Geodesy it is important to know that all kernels that are derived from Green-type potentials, e.g. by applying some differential operators, and taking the limit to the boundary, have this property. That holds in particular for the examples in Section 1.

The final result is that the integral (1) is defined as finite part integral provided that the density function is smooth enough, i.e. $\tilde{\mu}$ should be at least of class $C^{s-2, \alpha}(S)$. We write

$$\begin{aligned} \oint_S \tilde{K}(x, y-x) \tilde{\mu}(y) dS(y) &= \text{f.p.} \int_{S \setminus S_\varepsilon} \tilde{K}(x, y-x) \tilde{\mu}(y) dS(y) \\ &= \underbrace{\int_T (\text{ kernels of order } O(\rho^\sigma), \sigma \geq -1) dT}_{\text{absolutely integrable, cf. (13)}} + \underbrace{\sum_{i=0}^{s-2} \text{f.p.} F_i(\varepsilon)}_{\text{regular line integrals over boundary of } T} \end{aligned} \quad (37)$$

If the conditions (24) and (25) are fulfilled the Cauchy principal value exists, and we write

$$\begin{aligned} \oint_S \tilde{K}(x, y-x) \tilde{\mu}(y) dS(y) &= \lim_{\varepsilon \rightarrow 0} \int_{S \setminus S_\varepsilon} \tilde{K}(x, y-x) \tilde{\mu}(y) dS(y) \\ &= \underbrace{\int_T (\text{ kernels of order } O(\rho^\sigma), \sigma \geq -1) dT}_{\text{absolutely integrable, cf. (13)}} + \underbrace{\sum_{i=0}^{s-2} \text{p.v.} F_i(\varepsilon)}_{\text{regular line integrals over the boundary of } T} \end{aligned} \quad (38)$$

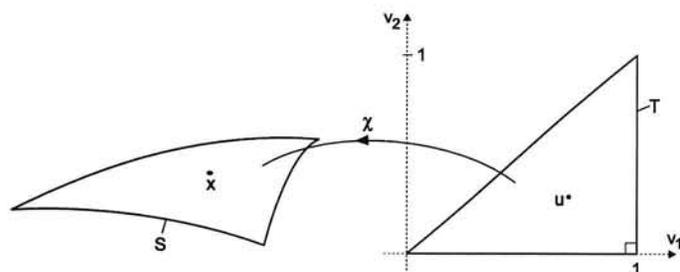


Fig. 6. Geometry of the numerical test

The regularization, which makes use of the definition of the finite part of a singular integral and of the Cauchy principal value, respectively, results in absolutely integrable integrals over the domain T and the sum of $s - 1$ regular integrals over the boundary of T .

4 Numerical integration

In order to test the formulas derived in Section 3, we consider a curved triangle S , which is the image of the standard triangle T under the quadratic mapping

$$\chi(v) := \begin{pmatrix} v_1 + v_2 + v_1^2 - 2v_2^2 + 4v_1v_2 \\ -v_1 + v_2 + v_1^2 - v_2^2 + 10v_1v_2 \\ v_1 + v_2 + v_1^2 - 2v_1v_2 \end{pmatrix}. \quad (39)$$

The density function is assumed to be linear in local coordinates, i.e.

$$\tilde{\mu}(y) = \tilde{\mu}(\chi(v)) = 1 + v_1 + v_2, \quad (40)$$

and the computation point is located at $u = \begin{pmatrix} 0.5 \\ 0.25 \end{pmatrix}$ (cf. Figure 6). As kernel functions we use the kernel of the double layer potential ($s=1$)

$$K_d(x, y) = D_{n(y)} \left(\frac{1}{|y-x|} \right) = \frac{\langle n(y), y-x \rangle}{|y-x|^3}, \quad (41)$$

the oblique derivative of the single layer kernel ($s=2$)

$$K_{os}(x, y) = D_{\tau(x)} \left(\frac{1}{|y-x|} \right) = \frac{\langle \tau(x), y-x \rangle}{|y-x|^3}, \quad (42)$$

and the oblique derivative of the double layer kernel ($s=3$)

$$K_{od}(x, y) = D_{\tau(x)} D_{n(y)} \left(\frac{1}{|y-x|} \right) = \frac{\langle n(y), \tau(x) \rangle}{|y-x|^3} - 3 \frac{\langle n(y), y-x \rangle \langle y-x, \tau(x) \rangle}{|y-x|^5}. \quad (43)$$

The following questions will be addressed:

1. Can finite part integrals be computed with the same accuracy as absolutely integrable integrals or do numerical instabilities occur?

2. Does the computation of finite part integrals require more nodes in order to achieve the same accuracy as absolutely integrable integrals or are many more nodes necessary?
3. Are the numerical costs in terms of function evaluations and number of arithmetic operations comparable or is the numerical effort per node higher?

As already has been shown (cf. (37),(38)), the regularization of the finite part integrals yields integrals with kernel functions of the order $O(\rho^\sigma)$, $\sigma \geq -1$, which are absolutely integrable, and $s-1$ regular integrals over the boundary of the integration domain. If $\sigma \geq 0$ we can directly make use of Gauss-Legendre product formulas because the kernels are analytic. If $\sigma = -1$, the integrals are weakly singular. Although they are still absolutely integrable, the numerical integration requires some care due to the singularity at $\rho = 0$. Let us briefly summarize how they can be computed efficiently; for more details see Klees (1996). The integrals have the form (cf. (13))

$$I = \int_T f(u, v) dT(v), \quad f(u, v) = O(\rho^{-1}). \quad (44)$$

We obtain

$$\begin{aligned} I &= \int_0^{2\pi} \int_0^{R(\varphi)} \underbrace{(f \rho)}_{=:k(\rho, \varphi)=O(1), \rho \rightarrow 0} d\rho d\varphi = \sum_{n=1}^3 \int_{\varphi_n}^{\varphi_{n+1}} \int_0^{R_n(\varphi)} k(\rho, \varphi) d\rho d\varphi \\ &= \sum_{n=1}^3 \int_{\varphi_n}^{\varphi_{n+1}} \int_0^1 \underbrace{k(\rho(\xi), \varphi) R_n(\varphi)}_{\text{oszillating in } \varphi} d\xi d\varphi = \sum_{n=1}^3 \int_{t_n(\varphi_n)}^{t_n(\varphi_{n+1})} \int_0^1 \underbrace{k(\rho(\xi), \varphi(t_n))}_{\text{analytic in } \xi, t_n} d\xi dt_n. \end{aligned} \quad (45)$$

The regularization has been performed using the Jacobian of the transformation of cartesian coordinates into polar coordinates. The new kernel $k(\rho, \varphi)$ has order $O(1)$, $\rho \rightarrow 0$. The integral over T has been split-up into three integrals, since the boundary $\rho = R(u, \varphi)$ of T is not smooth as φ runs from 0 to 2π . Each of the three integrals is smooth w.r.t. φ , because the boundary $R_n(u, \varphi)$, $n = 1, 2, 3$, is smooth. The transformation $\rho \rightarrow \xi : [0, R_n(u, \varphi)] \rightarrow [0, 1]$ yields a new kernel $k(\rho, \varphi) R_n(u, \varphi)$, which oszillates if the computation point u is near the boundary $\rho = R_n(u, \varphi)$. For instance, if $n = 2$ we have (cf. Figure 7)

$$R_2(u, \varphi) = \frac{h_2}{\cos \varphi}. \quad (46)$$

For small h_2 the integration bounds tend to $\pm \frac{\pi}{2}$, and $R_2(u, \varphi)$ has strong gradients. The practical implication is that then the integration over φ requires many nodes. In order to avoid this we can smooth the kernel by applying for each n a parameter transformation $\varphi \rightarrow t_n : [\varphi_n, \varphi_{n+1}] \rightarrow [t_n(\varphi_n), t_n(\varphi_{n+1})]$, defined by

$$dt_n = R_n(u, \varphi) d\varphi, \quad n = 1, 2, 3. \quad (47)$$

Then, the new kernel $k(\rho(\xi), \varphi(t_n))$ is analytic w.r.t. the new variables t_n and ξ , and Gauss-Legendre product formulas can be applied to compute the integral efficiently.

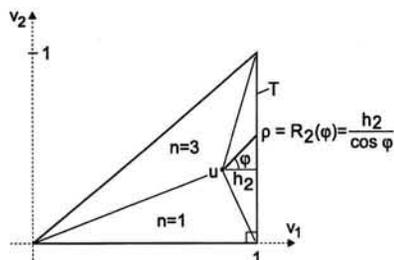


Fig. 7. Regular line integrals over the boundary of T

Table 1. Number of nodes and relative integration error for various kernel functions

	number of nodes, relative integration error							
K_d	12	1.0(-3)	27	1.6(-6)	48	3.4(-6)	108	2.0(-9)
K_{os}	18	6.8(-3)	36	9.5(-5)	60	3.3(-6)	126	4.0(-10)
K_{od}	18	4.1(-3)	36	4.2(-5)	60	6.1(-6)	126	7.2(-9)

The regular line integrals over the boundary of T have the form

$$I = \int_0^{2\pi} g_i(u, \varphi) \left\{ \begin{array}{ll} \ln R(u, \varphi) - \ln d_0(u, \varphi), & i = s - 2 \\ -\frac{1}{m} [R^{-m}(u, \varphi) - d_{m,m}(u, \varphi)], & 0 \leq i < s - 2 \end{array} \right\} d\varphi, \quad (48)$$

with $m = s - 2 - i$ and $0 < m \leq s - 2$. The integrand is not smooth w.r.t. φ for the same reason as mentioned before. Again, we have to split-up the line integral into three integrals over $\rho = R_n(u, \varphi)$, $n = 1, 2, 3$, and have to perform parameter transformations $\varphi \rightarrow t_n$, $n = 1, 2, 3$, in order to get an analytic integrand. Then, Gauss-Legendre quadrature formulas are the best choice in order to get high accuracies with a minimum number of nodes.

The results of the test calculations are shown in Table 1 and Table 2. Table 1 gives answer to the first two questions: there is no difference between finite part integrals and absolutely integrable integrals in terms of accuracy and number of nodes. This implies that finite part integrals can be calculated as accurately as absolutely integrable integrals, and do not require more nodes to achieve the same accuracy. In particular, the regularization does not cause any instabilities. Striking is that only very few nodes are needed in order to make the integration error small.

Table 2 gives answer to the third question: the numerical costs in terms of number of arithmetic and function evaluations is much higher for finite part integrals than for absolutely integrable integrals. This is caused by the regularization of the singular integrals, which yields kernel functions that are more elaborate to compute than the kernel of the singular integral. Therefore, the higher the order of singularity is the higher the numerical effort per node.

Table 2. Number of arithmetic operations and number of function evaluations for various kernel functions

	K_d	K_{os}	K_{od}
number of arithm. operations	$489 M^2$	$392 M^2 + 18 M$	$3484 M^2 + 52 M$
	4401	3582	31512
	14%	11%	100%
number of function evaluations	$10 M^2$	$15 M^2 + 11 M$	$118 M^2 + 23 M$
	90	168	1131
	8%	15%	100%

M =number of nodes of the Gauss-Legendre quadrature formula; relative integration error $\approx 10^{-4}$ if $M = 3$.

5 Summary

We discussed the concept of strongly singular and hypersingular surface integrals with kernels possessing a point singularity of order ≥ 2 in the interior of the domain of integration. We have shown how to transform them into absolutely integrable integrals over the parameter domain and regular one-dimensional integrals over the boundary of that domain. The choice of the neighborhood of the computation point in the parameter domain is in general not arbitrary but depends on the choice of the neighborhood on the surface. Only if the kernel function has some special properties the finite part is invariant w.r.t. the choice of the neighborhood in the parameter domain. Then, the most convenient neighborhood, the circle centered at the computation point, can be used.

The regularized two-dimensional integrals that have weak point singularities can be computed efficiently using Gauss-Legendre product formulas after some special parameter transformations are applied. For a given number of nodes, strongly singular and hypersingular integrals can be computed with the same accuracy as absolutely integrable integrals. The numerical costs per node, however, do strongly depend on the order of singularity. Thus, hypersingular integrals are much more costly to calculate than strongly singular integrals, and the computation of the latter takes more time than the computation of weakly singular integrals. This is due to the regularization, which yields complicated expressions for the regularized kernels, which are more elaborate to evaluate than the original kernels.

The one-dimensional line integrals can be calculated without any problems using Gauss-Legendre quadrature formulas. Only if the computation point is near the boundary, additional parameter transformations have to be applied before Gauss-Legendre quadrature formulas are used. The same holds for the regularized two-dimensional integrals.

Acknowledgement The figures have been made by A. Smits from the Faculty of Civil Engineering and Geo Sciences. His support is gratefully acknowledged.

References

- Bian, S. and Sum, H. (1994). The expression of common singular integrals in physical geodesy. *manuscripta geodaetica*, **19**, 62–69.
- Bosch, W. (1977). Geschlossene integration von potentialkernen beim modell der einfachen schicht mit stückweise ebenem rand. Mitteilungen aus dem Institut für theoretische Geodäsie 49, Universität Bonn.
- Guiggiani, M. (1991). The evaluation of cauchy principal value integrals in the boundary element method - a review. *Math Comput. Modelling*, **15**, 175–184.
- Guiggiani, M., Krishnasamy, G., Rudolphi, T., and Rizzo, F. (1992). A general algorithm for the numerical solution of hypersingular boundary integral equations. *Transaction of the ASME*, **59**, 604–614.
- Hofmann-Wellenhof, B. (1983). Representation of the gravitational potential by multipoles. Mitteilungen der geodätischen institute folge 47, Technical University Graz.
- Kieser, R. (1991). *Über einseitige Sprungrelationen und hypersinguläre Operatoren in der Methode der Randelemente*. Ph.D. thesis, University of Stuttgart.
- Kieser, R., Schwab, C., and Wendland, W. (1992). Numerical evaluation of singular and finite-part integrals on curved surfaces using symbolic manipulation. *Computing*, **49**, 279–301.
- Klees, R. (1992). *Lösung des fixen geodätischen Randwertproblems mit Hilfe der Randelementmethode*. Ph.D. thesis, Deutsche Geodätische Kommission, Reihe C, No. 382, München.
- Klees, R. (1996). Numerical calculation of weakly singular surface integrals. *Journal of Geodesy*, **70**, 781–797.
- Meissl, P. (1971). Preparations for the numerical evaluation of second order molodensky-type formulas. Reports of the department of geodetic science no. 163, Ohio State University.
- Schwab, C. and Wendland, W. (1992a). Kernel properties and representations of boundary integral operators. *Mathematische Nachrichten*, **156**, 187–218.
- Schwab, C. and Wendland, W. (1992b). On numerical cubatures of singular surface integrals in boundary element methods. *Numerische Mathematik*, **62**, 343–369.
- Shaofeng, B. and Xurong, D. (1991). On the singular integration in physical geodesy. *manuscripta geodaetica*, **16**, 283–287.
- Vijayakumar, S. and Cormack, D. (1988). An invariant imbedding method for singular integral evaluation on finite domains. *SIAM J. Appl. Math.*, **48**, 1335–1349.

1870

1871

1872

1873

1874

1875

1876

1877

1878

1879

1880

1881

1882

1883

1884

1885

1886

1887

1888

1889

1890

1870

1871

1872

1873

1874

1875

1876

1877

1878

1879

1880

1881

1882

1883

1884

1885

1886

1887

1888

1889

1890

The Fiction of the Geoid*

Roland Klees and Martin van Gelderen

Abstract

We discuss the role of the static and time-varying geoid in sea-level studies by means of three examples, i.e. (a) mean ocean circulation, (b) post-glacial rebound, and (c) vertical datum connection. First, the contribution of the geoid to these items is addressed. Then, the requirements in terms of accuracy and resolution are discussed and compared with the current knowledge about the static and time-varying geoid. Thereafter, we discuss the main problems of current geoid determination in terms of data, models, and numerics. Finally, we show the impact of dedicated satellite gravity mapping missions on the quality of the static and time-varying geoid in terms of accuracy and spatial and temporal resolution, and discuss the implications of an improved geoid for sea-level change studies. Currently, the contribution of the geoid to sea-level studies is rather weak for various reasons such as data distribution, data quality and data handling, inconsistency in the mathematical models used, and numerical and conceptual problems. All make the geoid a fiction when talking about accuracies of 10^{-8} and higher over various spatial and temporal scales as needed in sea-level studies. This will change dramatically, however, if a dedicated gravity field mission will be launched provided that at the same time adequate functional models and numerical techniques are at our disposal.

1 Introduction

It is not a secret that the geophysical mechanisms of sea level change are not fully understood and that current estimates of future sea level change are not sufficiently accurate and reliable. However, most of the processes relevant for sea level studies involve redistribution of mass, and, therefore, are likely to result in changes in the geoid. This sensitivity of the geoid w.r.t. mass redistribution gives rise to the question whether we can learn more about the mechanism of mass exchange between ice and ocean and the reaction of the deformable Earth to these forces by studying the inverse problem, i.e. to deduce from measured changes in the geoid information about the underlying processes. In order to answer this question, we have to study the sensitivity of the geoid w.r.t. mass redistribution and the quality of current and future geoid models. In addition, the problem of separation of competing sources of geoid variability has to be addressed, because there are other processes, not related to climate, which may cause geoid variations. Besides the time-varying geoid also the static geoid plays a significant role

* Pres. at the *Staring Symposium*, October 21, 1997, Royal Academy of Sciences, Amsterdam, the Netherlands

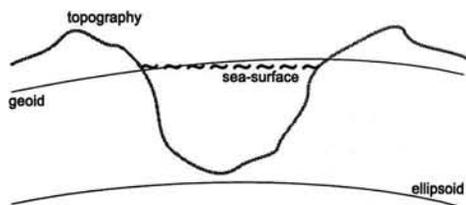


Fig. 1. The geoid and the mean sea surface

in sea level studies; for instance, by providing a reference surface in order to model relevant dynamic processes such as mean ocean circulation, in order to connect different continental height systems, and in order to describe absolute sea level changes in a unified world height system.

The goal of the paper is to discuss by means of some examples the current and future role of the static and time-varying geoid in sea-level studies. In Section 2 we will make some remarks about the geoid. In Section 3 we will illustrate the potential implications of the geoid for sea-level studies by means of three examples, i.e. ocean circulation, post-glacial rebound, and vertical datum connection. Section 4 is devoted to the current knowledge about the geoid and the main problems of current geoid determination. In Section 6, finally, we will discuss the implication of future geoid models for sea-level studies.

2 The geoid

The gravity field of the Earth is the net result of the Newtonian gravitational attraction of the Earth masses and the Earth rotation. To a first approximation, i.e. to better than 10^{-5} , the Earth has an ellipsoidal shape and gravity field. Therefore, the shape and gravity field of the Earth is commonly expressed as departure from a well-defined ellipsoidal reference field. One possibility to quantify this departure is by means of the vertical separation between corresponding equipotential surfaces of the real and the reference gravity field. For the real field this is the geoid, defined as the equipotential surface of the real gravity field that closely corresponds with the mean sea level; for the reference field it is the surface of the ellipsoid of revolution which is an equipotential surface of the reference gravity field. The vertical separation between geoid and level ellipsoid is called "geoid height" (see Figure 1). The time-averaged or mean sea surface, to which elevations on land are referred, is vertically displaced from the geoid since the fluids in the oceans are in motion w.r.t. the solid Earth. These departures are on the order of 1-2 m, about two orders of magnitude less than the geoid heights. They are referred to as "mean dynamic sea surface topography". In absence of wind-driven currents, heat-driven and salt-driven circulation, and river discharges, the mean dynamic sea surface topography would be equal to zero, meaning that mean sea surface and geoid would essentially coincide. From space we can measure the height of the mean sea surface above the reference ellipsoid, which is the sum of the geoid height and the mean dynamic sea surface topography.

The current state-of-the-art geoid model, the EGM96, is shown in Figure 2. The geoid heights approach maximum values of about 100 m, the global RMS is about 42 m. They reflect in fact the dynamics of the Earth caused by a variety of processes on a wide range of spatial

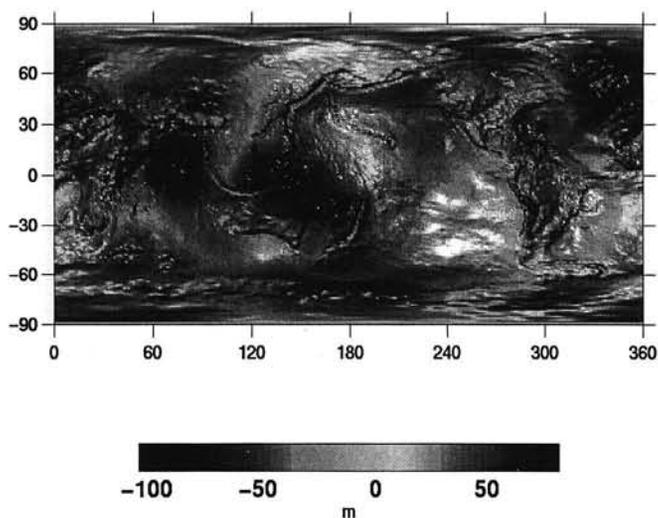


Fig. 2. The EGM96 geoid w.r.t. GRS80

and temporal scales ranging from kilometers to worldwide and from hours to million of years. Dominant are irregularities in the solid Earth caused by convective processes that deform the Earth on time scales of thousands to millions of years. The time-varying part of the geoid takes less than 1% of the total geoid signal over human life-times (cf. National Research Council (1997)).

3 Geoid and global sea level change

The role of the static and time-varying geoid in sea level studies is determined by the sensitivity of the geoid w.r.t. mass distribution and mass transfer among the solid Earth, the hydrosphere, cryosphere, and atmosphere, respectively. That is

1. to discriminate among causes of variation in sea level, e.g. between ocean thermal expansion and mass inflow as sources of a sea level change;
2. to improve our understanding of processes causing sea level change, e.g. by studying large scale ocean circulation or by determining changes in the mass distribution of polar ice;
3. to provide constraints to models of geophysical processes affecting global sea level such as post-glacial rebound;
4. to serve as a reference surface w.r.t. which absolute changes are measured;
5. to provide high accurate orbits of altimeter satellites, which map the sea surface on a global scale with high spatial and temporal resolution.

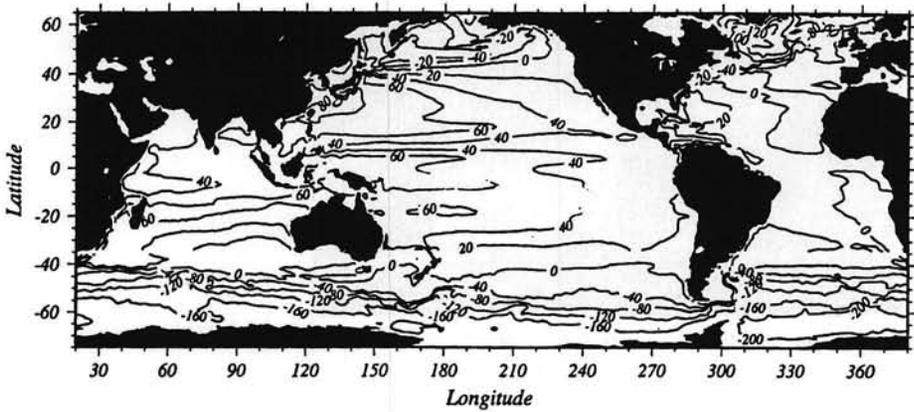


Fig. 3. The Dynamic Sea Surface Topography (from ESA (1997))

Let us make this clear by means of three examples.

The main problem in interpreting global sea level is the separation between the total mass of water in the ocean, the mean temperature of water, and the ocean circulation (e.g. Wunsch (1993)). The geoid, in conjunction with altimeter measurements, determines the time-averaged or mean dynamic sea surface topography. An example is shown in Figure 3. Particularly important in the context of ocean circulation are the basin-wide currents with slopes of about 10^{-8} and scales between 1000 km and 3000 km, the Antarctic Circumpolar Current with typical slopes of about 10^{-6} and scales of 500-1000 km, and the Western Boundary Currents such as the Gulf Stream and the Kuroshio with slopes on the order of 10^{-5} , and typical scales between 50-100 km (cf. National Research Council (1997); ESA (1997)). Since the slope of the mean dynamic sea surface topography is proportional to the mean surface geostrophic velocity, the geostrophic surface currents can be determined. These, in turn, can be used with in-situ measurements of temperature and salinity, in inverse calculations of the deep ocean circulation, i.e. of the transport of sea water, heat, and salt, which is a key factor in regulating the Earth's climate on decadal and longer time scales. Therefore, the role of the geoid has to be seen in the context of interpretation, understanding, and prediction, as opposed to the monitoring, of global sea level change.

A nice example how the geoid can provide observational constraints to models of geophysical processes affecting global sea level is post-glacial rebound. Post-glacial rebound is related to sea-level change in two respects:

Firstly, post-glacial rebound induces a sea-level signal because it affects the topography of the Earth and the geoid. For instance, the top panel of Figure 4 shows the secular change in the geoid due to post-glacial rebound over North America (National Research Council (1997)). The maximum amplitude is about 2.4 mm/yr. The topographic signal can even be larger, up to 10 mm/yr (National Research Council (1997)). These signals must be subtracted from the altimeter and tide-gauge measurements in order to quantify the climatic contribution to sea-level change, e.g. due to thermal expansion of ocean water and the melting of glaciers and ice

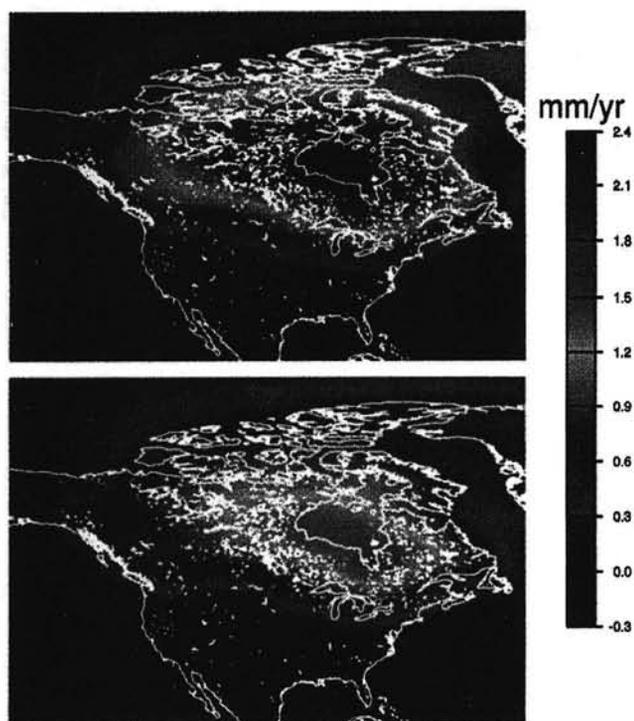


Fig. 4. Post glacial rebound from two lower-mantle viscosity models (from National Research Council (1997))

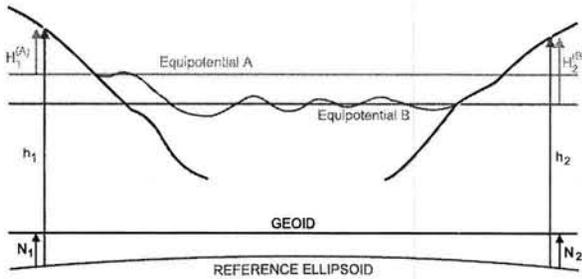


Fig. 5. The principle of vertical datum connection

sheets. Therefore, accurate knowledge about the post-glacial rebound signal in the geoid and the topography is needed.

Secondly, post-glacial rebound is indispensable to make reliable predictions of future sea-level change. Since it is mainly controlled by the mantle viscosity, most of the present-day post-glacial-rebound studies try to provide constraints on the viscosity profile of the mantle. Very important in that respect are independent observational constraints for the *lower*-mantle viscosity in order to validate the geophysical models. They are hard to get due to the huge masses needed, but can be provided by study of the secular change in the geoid due to e.g. the melting of the Wisconsin ice sheet over North America. The reason is that different lower-mantle viscosity values result in different amplitudes of secular geoid variations, but remains the geoid pattern almost unchanged (cf. Figure 4). For instance, the geoid signal shown in the bottom panel is based on a lower-mantle viscosity 5 times larger than in the top panel. This yields a maximum geoid signal of only 1.5 mm/yr compared to the 2.4 mm/yr for the lower viscosity value. Therefore, measuring temporal variations in the geoid would allow to resolve differences between competing viscosity models.

A last example concerning the role of the geoid in global sea level studies is the connection of height systems. Usually, the origin of a continental height system has been chosen by assigning a value to a reference marker, a bench mark (cf. Figure 5). The value chosen was based on the observation of the mean sea level for a given time interval at a time one believed that mean sea level and geoid are identical. Since in fact mean sea surface and geoid deviates up to 2 m and due to non-homogeneous sea level changes, land movements, and local currents, each region has in fact its own vertical datum. One implication is that tide gauge data referring to different height systems cannot be used in global studies of absolute mean sea level. The geoid would provide a physically meaningful connection between the various height systems, which in turn would allow to transform all tide-gauge time series into a unified world height system.

4 Current knowledge of the geoid

Let us now have a look at the current knowledge about the geoid and whether this knowledge is sufficient for sea level studies in terms of accuracy and resolution.

Figure 6 illustrates the development of the knowledge about the geoid in terms of resolution and accuracy. The first global geoid models, which date back to the thirties, were based on terrestrial gravity data and limited to the very low degrees, thus representing wavelengths of, say, 5000 to 10000 km. The advent of satellite altimetry in 1975 revolutionized the measurements

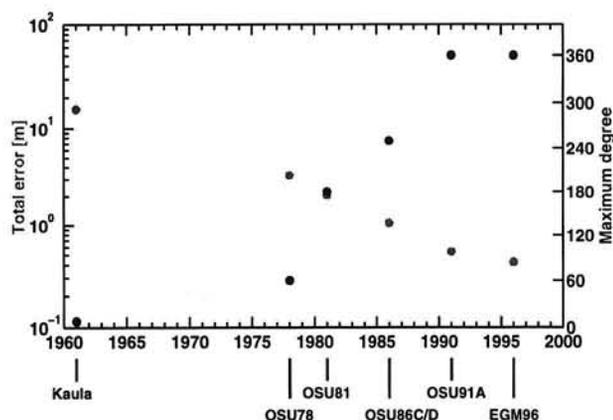


Fig. 6. Development of the knowledge of the gravity field. The light dots refer to the total error (left scale) the darker dots to the maximum degree and order of the of spherical harmonic expansion (right scale)

of much shorter wavelengths of the marine geoid resulting in a geoid model of degree 52 in 1978, which corresponds to wavelengths of about 800 km. Progress in terrestrial gravimetry, satellite-tracking techniques, satellite altimetry, and gravity solution techniques, together with an increasing number of satellites for which data is available, have led to the development of many gravity models with increased resolution over the last 20 years. The most recent global gravity field, the EGM96, represents all wavelengths up to degree 360, corresponding to a spatial resolution of about 55 km at the equator.

The same holds true for the development of the accuracy. Starting from some 20 meter total RMS in the geoidal heights in the early sixties, which amounts to about 50% of the total RMS signal, we have now approached a level of about 40 cm global RMS for the most recent model, due to more data, more accurate data and more satellites. However, the quality is highly non-homogeneous, and there are still areas almost not covered with gravity data, where the geoid error can reach several meters.

What does this imply for the role of the geoid in sea level studies, e.g. for determining the geostrophic surface currents? Figure 7 shows the geoid slope errors versus the spatial resolution for the current state-of-the-art geoid model. This is compared with the mean sea surface slope errors as derived from satellite altimeter measurements, and with the slopes and scales of the currents we want to recover, namely the basin-wide currents (BASIN), the Antarctic Circumpolar Current (ACC), and the Western Boundary Currents (WBC). Obviously, altimetry fits the requirements except for the very long wavelengths. The geoid slope error, however, exceeds the altimeter error at resolutions shorter than about 3000 km. Therefore, the present ability to resolve the most important currents is limited by the geoid slope errors at the corresponding scales.

The next example was post-glacial rebound. Here we expected that the geoid can provide unambiguously observational constraints to the viscosity profile of the mantle. However, present-day observed secular changes in the geoid are limited to a few very long-wavelength

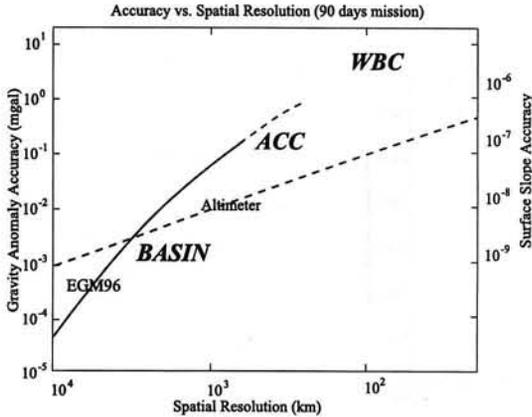


Fig. 7. Accuracy vs. resolution of the EGM96 gravity model (EGM96); refer to the left hand scale. The surface-slope scale is shown on the right-hand scale with the altimetry slope error (assuming 10mm uncertainty in altimetry height differences), the approximate slope of basin-wide currents (BASIN), the Antarctic Circumpolar Current (ACC) and the Western Boundary Currents (WBC). (From National Research Council (1997))

components, which are derived from satellite laser tracking. The Figure 8, e.g. shows the result of five years observation of changes in the difference between the Earth's polar and equatorial moments of inertia. The observed variations show strong annual and inter-annual components. They do certainly contain a post-glacial-rebound signal, but the dominant features are the cumulative effect of a number of processes such as atmospheric mass redistribution, long-period non-equilibrium ocean tides, continental water storage, snow cover, and ice sheet volume changes. A separation of the post-glacial rebound signal from the other effects is hardly possible with the currently available information. In order to do that information about secular changes in the geoid on much shorter wavelengths is indispensable.

The last example shown referred to the problem of vertical datum connection. Figure 9 shows the estimated RMS for the vertical datum connection between various height systems on the world. The left panel uses a geoid model which has been derived from purely satellite tracking data, the right panel uses a model including terrestrial gravity data. Depending on the regions to be connected, we estimated errors up to 80 cm when using a satellite only geoid and about 20 cm if terrestrial gravity data is taken into account. Although some other error sources contribute to the total budget, the geoid error is by far dominant. Therefore, a datum connection at a level of, say, some centimeters, requires a more accurate geoid over various wavelengths depending on the spatial distance between different height systems.

The vertical datum inconsistencies may have serious consequences, e.g. for the geostrophic velocity estimates derived from geoid slopes and mean sea surface slopes. To illustrate this we have calculated the effect on the geoid of the vertical datum difference between the Dutch and the British height system. The apparent geoid heights shown in Figure 10 result in wrong estimates of the geoid slope, which in turn imply erroneous geostrophic velocity field estimates over that area.

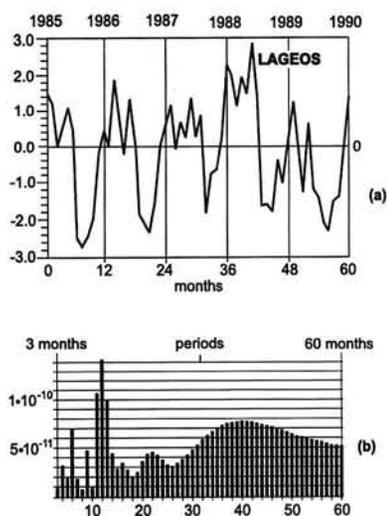


Fig. 8. Change of the difference between the earth's polar and equatorial moments of inertia. Upper panel shows the time series as found from Lageos tracking data; the bottom panel its spectrum (from Gegout and Cazenave (1993)).

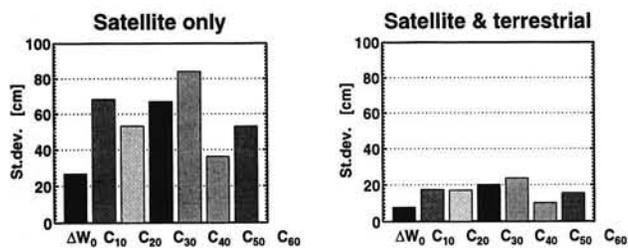


Fig. 9. Vertical datum connection: current situation (from Onselen (1997))

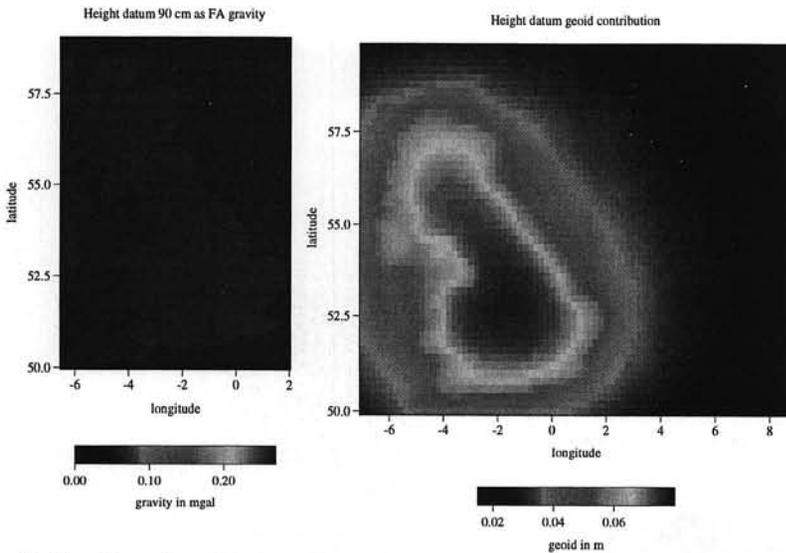


Fig. 10. The effect of vertical datum inconsistencies on geoid heights

5 State-of-the-art geoid - the problems

Obviously the current knowledge about both the static geoid and the time-varying geoid is too weak for many global sea level studies. One reason is the data distribution and data quality, although state-of-the-art geoid models, incorporate virtually all available data. That are (a) satellite tracking data, including optical data, radio Doppler and radio interferometry observations, satellite laser ranging, and microwave tracking data from GPS, DORIS and PRARE, (b) about twenty years of satellite altimetry, (c) measured surface gravity data (land, sea, air) and geophysically predicted gravity data, and (d) digital elevation models of the Earth's topography. Each of these data sets contributes to the geoid model in a different way and has its own deficiencies and limitations.

Satellite tracking data, for instance, provides a global data coverage, and a rather homogeneous data quality, but can only determine the long-wavelength components of the geoid, i.e. wavelengths of a few thousand kilometers and longer. One reason is the well-known attenuation of the gravitational attraction with increasing distance to the mass. To explain this effect let us have a look at Figure 11. The bottom panel shows the surface gravity data, which is closely related to the Earth's topography, thus containing many fine, i.e. short-wavelength structures. The middle panel shows what is seen in satellite tracking data to satellites at an altitude of 450 km. All the fine structures of the gravity field have been smoothed out and only the dominant large scale features are still visible. The practical implication of the attenuation effect is that from noisy satellite tracking data the short-wavelength features in the geoid cannot be recovered.

Surface gravity data is in principle capable of resolving all wavelength features of the geoid provided that a uniform dense global coverage of high-quality is available. However, surface gravimetry is a time-consuming and expensive measurement technique and current data sets are derived from several thousand different sources collected over decades with different instru-

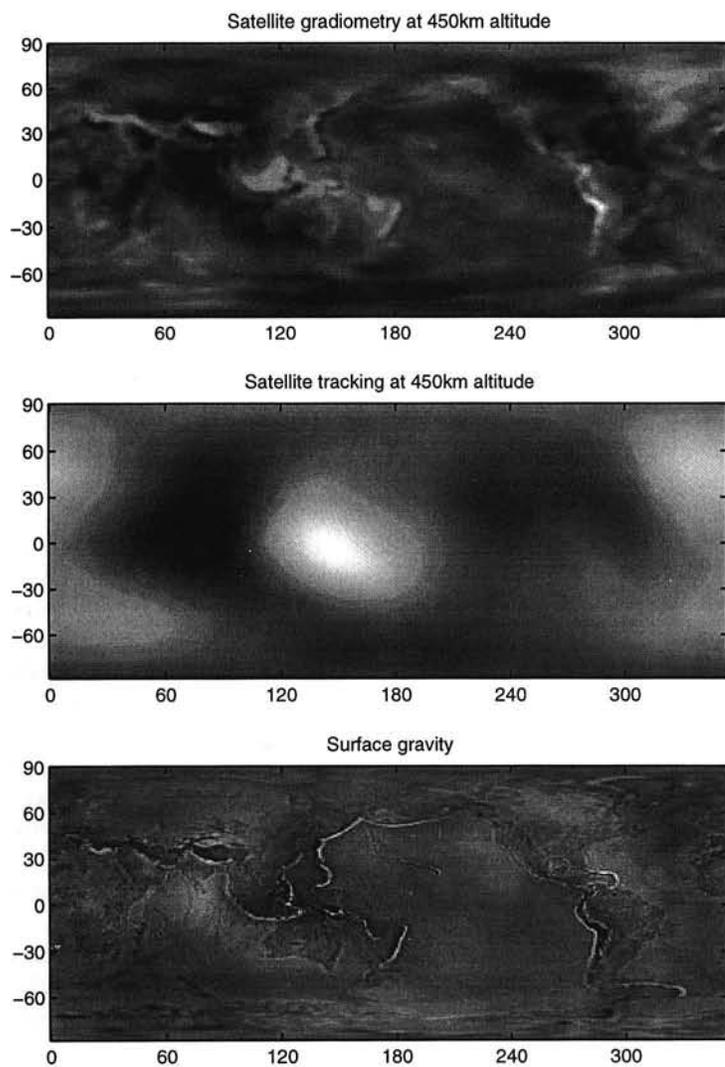


Fig. 11. Signal of second vertical derivative of the potential (top), tracking at satellite altitude (middle) and surface gravity data (bottom)

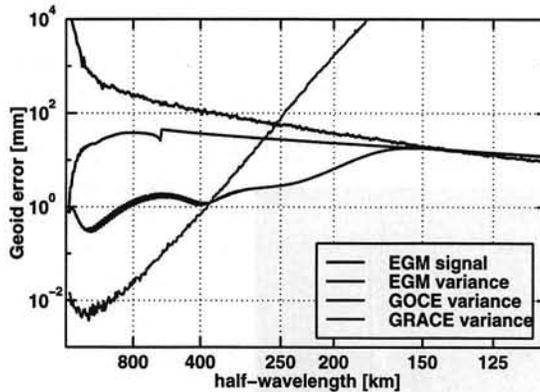


Fig. 12. The geoid errors by degree for several models

ments. Although the data holdings have improved dramatically over the last years, many lakes, high mountain areas, shallow water areas, and polar regions are almost void of gravity measurements. In addition, the accuracy and density of gravity data vary substantially with geographic region and the gravity data sets are contaminated by systematic errors.

Satellite altimetry provides an unsurpassed mapping of the (mean) sea surface in terms of accuracy and resolution. However, altimetric measurements are confined to the ocean areas and the satellite's inclination leaves high latitude areas uncovered. Moreover, the conversion of altimetric measurements into gravity data requires the time-averaged dynamic sea surface topography to be known. Both aspects weaken the quality of geoid heights and geoid slopes derived from altimeter measurements.

The problems related to data distribution, coverage and quality can only be solved by means of satellite techniques, provided that the attenuation problem can be solved, as well. The latter, however, can be addressed by choosing a new observation type at satellite altitude. For instance, the top panel in Figure 11 shows what can be recovered at satellite altitude when classical satellite tracking data is replaced by measurements of the second derivatives of the geopotential. Obviously, the sensitivity w.r.t. the short-wavelength features improves dramatically. Significant improvements are also possible if range rates are measured between two satellites.

Therefore, a number of studies have been done, and are still in progress, for a dedicated gravity field mission, which will improve all wavelengths of the geoid down to wavelengths of several hundred kilometers. Two mission designs are considered. The so-called satellite-to-satellite tracking (SST) utilizes differential tracking of two satellites and thereby measures orbital perturbations. The so-called satellite gravity gradiometry (SGG) measures the differences in acceleration of two masses within the same spacecraft. Examples are ESA's GOCE mission, which utilizes a combination of SGG and high-low SST, and the U.S. GRACE mission, which utilizes low-low SST only.

In order to get an idea what really will improve, we can compare the geoid height error by degree for the GOCE and the GRACE mission with the current state-of-the-art geoid model EGM96 (Figure 12). Clearly, the improvement is dramatic down to half-wavelengths of about 200 km for GOCE and 400 km for GRACE.

But data distribution and data quality are not the only problems in current geoid determina-

tion; modeling and numerical aspects are important, as well.

1. For instance, the standard approach to the solution of the geoid is the discrete weighted least-squares technique which is used in order to get a best linear estimate w.r.t. a hybrid norm. The functional relationship between observations and gravity field parameters defines the functional model which is a mathematical description of the physical reality. Together with the gravity field parameters, many other parameters will appear in the functional model as well, e.g. parameters describing station coordinates, satellite orbits, non-conservative forces, atmospheric path delay, tides etc. The quality of the functional model depends on how well it represents reality, and it is by far not a trivial problem to set up a consistent functional model in the sense that the predefined goals are met w.r.t. accuracy and resolution.
2. The functional model will in general be highly non-linear. Therefore, a linearization is performed based on appropriate reference models. The neglected higher-order terms are attributed to the model noise. The set-up of an iteration process is obvious, but by far not trivial, and suitable iteration procedures have not been designed and applied yet.
3. Although millions of discrete observations will enter the model, the system is always highly under-determined due to the structure of the gravity field. An approximate solution is sought in a finite dimensional subspace by limiting to a finite number of gravity field parameters to be estimated. The resulting error is again attributed to the model error, assuming that this is consistent with the requirements.
4. Not all parameters are well determined by the observations resulting in unstable normal equations which requires some regularization. This causes some biases in the weakly determined parameters, which strongly depends on the regularization. In addition, the finite dimension of the solution space and the neglect of higher order terms directly contribute to the bias.

Other problems are related to the feasibility of a discrete least-squares approach.

1. For, e.g. to resolve all wavelength down to 110 km, about 130000 gravity field parameters are required. In addition, thousands of nuisance parameters have to be included into the functional model.
2. Limited computational capabilities make a proper discrete least-squares approach not feasible at the time being, but require alternative strategies and simplifications. Consequently, current geoid models, even when based on the same data sets, show global RMS differences of some decimeters depending on the followed strategy, data handling, and simplifications. For instance, non-global data and overlapping data are "removed" by applying approximation techniques; satellite altimeter data are not treated as direct tracking observations but are converted into gravity data; observations are not weighted individually according to their estimated error variances; existing correlations are not taken into account, and observations are interpolated on some regular grid to make them suited for fast numerical techniques.
3. Many heterogeneous observations enter the functional model. The stochastic model is incomplete; correlations are mostly neglected, and even realistic variance estimates are often not available. Moreover, numerous reductions have to be applied to the data, among them reductions to further limit the numerical effort.

4. Finally, the observations are contaminated by systematic errors, e.g. vertical datum inconsistencies, inconsistent data reductions, and instrumental biases.

The high accuracies needed in sea level studies make even a consistent conceptual definition of the geoid an elusive and by far not trivial procedure. The definition of the geoid as equipotential surface that closely corresponds with the mean sea level is sufficient when accuracies on the order of 10^{-6} are needed. The same holds true for other more fancy definitions such as the geoid as the surface of a homogeneous ocean under the influence of the Earth's gravity field or the geoid as the surface of a uniform and static ocean. However, none of them is the definition for use in sea level studies where accuracies of 10^{-8} and higher are needed depending on the scale. That is because they give not sufficient credit to the observables to be used in implementing the geoid, to the temporal variations in the geoid, and to the existence of a dynamic sea surface topography.

Many alternatives have been proposed in the past. Each identifies a different surface as the geoid. However, none of them is without problems, especially from an operational point of view. For instance, most definitions assume that the total tidal effect has been removed, including the permanent tides; an effect which amounts up to 3 decimeters in geoidal heights. This, however, is not possible because we do not know the correct values of the Love numbers. Even though we can get an internally consistent model by using the same Love numbers, a "non-tidal" geoid will remain a fiction. Moreover, any lack of global coverage with measurements results in a realization which is only approximate.

6 Potential of future geoid models

Let us assume that the conceptual problem has been solved, that we have a consistent functional model at our disposal, and that gravity data is available from the GOCE and the GRACE mission. What would be the implications of the new geoid models for sea level studies?

The first example of Section 3 referred to dynamic sea surface topography and mean ocean circulation. Figure 13 shows what will happen with the geoid slope errors compared with the quality of altimeter derived mean sea surface slopes and with the typical scales and slopes of the most important currents we want to recover. Clearly, GOCE and GRACE will allow BASIN and ACC scales to be resolved accurately. The geoid slope error will even become insignificant in the range between 300-3000 km half-wavelengths, which is completely reverse to the current situation. Only the WBC scales will hardly be resolved because the geoid slope uncertainty will have about the same order as the slope of the dynamic sea surface topography at these spatial scales.

The second example was related to post-glacial rebound and observational constraints on lower-mantle viscosity. Figure 14 indicates that the post-glacial rebound signal can now readily be seen in the geoid. That is because the geoid signal of post-glacial rebound exceeds the uncertainty levels of a SST mission like GRACE at wavelengths greater than 1500 km. Secondly, differences between competing viscosity values can be resolved. For instance, the geoid signal of the difference between two viscosity values, which differ by a factor of 5, exceeds the GRACE uncertainty level at wavelengths greater than, say, 3400 km.

The last example referred to the vertical datum connection. The left panel in Figure 15 shows the current situation, i.e. the total RMS error for the connection between various vertical datums. The right panel shows what really improves if the geoid derived from measurements of the GOCE mission is used. The total error reduces from some 20 centimeters to less than 5

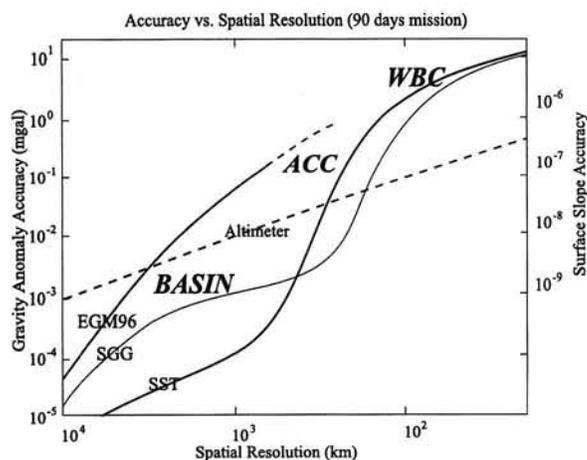


Fig. 13. Accuracy vs. resolution of the EGM96 gravity model (EGM96) and gravity fields to be obtained from a SGG or SST mission. The surface-slope scale is shown on the right-hand scale with the altimetry slope error (assuming 10mm uncertainty in altimetry height differences), the approximate slope of basin-wide currents (BASIN), the Antarctic Circumpolar Current (ACC) and the Western Boundary Currents (WBC). (From National Research Council (1997))

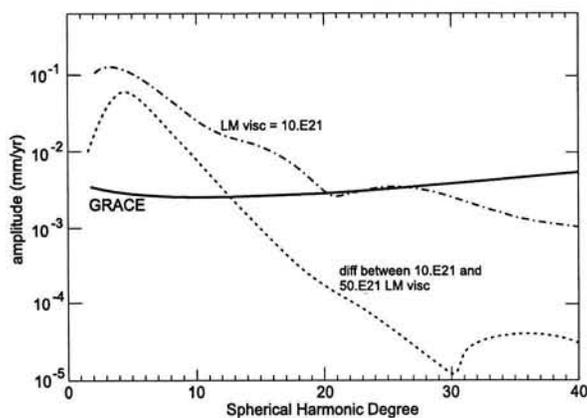


Fig. 14. The degree amplitudes for post-glacial rebound in North America for a lower-mantle viscosity of $10.E21$ Pa-sec and for the difference between results for viscosities of $10.E21$ and $50.E21$ Pa-sec, compared with degree variances of GRACE (from National Research Council (1997))

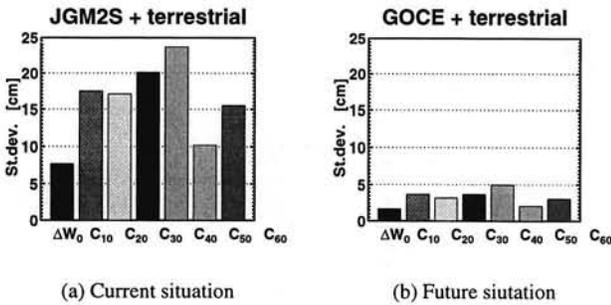


Fig. 15. Vertical datum connection quality with future GOCE data (from Onselen (1997))

centimeters. Moreover, when analyzing the total error budget we see that the geoid error will no longer be the dominating error source.

References

- ESA (1997). Assessment report of the gravity field and ocean circulation mission. , ESA SP-1196.
- Gegout, P. and Cazenave, A. (1993). Temporal variations of the earth's gravity field for 1985-1989 derived from Lageos. *Geophys. J. Int.*, **114**, 347-359.
- National Research Council (1997). *Satellite gravity and the geosphere*. National Academy Press, Washington, D.C.
- Onselen, K. V. (1997). Quality investigation of vertical datum connection. DEOS report 97.3, Delft Institute for Earth-Oriented Space Research (DEOS).
- Wunsch, C. (1993). Physics of the ocean circulation. In R. Rummel and F. Sanso, editors, *Satellite altimetry in geodesy and oceanography*, volume 50 of *Lecture Notes in Earth Sciences*, pages 10-99. Springer-Verlag, Berlin.

Natural Gas Extraction and its Induced Gravity Change

Martin van Gelderen, Roger Haagmans and Mirjam Bilker¹

¹ Finnish Geodetic Institute, Helsinki

Abstract

Above the Groningen gas field gravity observations are available over a 18 year time span. The old gravity data was reanalyzed and a new campaign executed. The observed gravity changes (from 4 campaigns) were compared with the gravity effect computed from the reservoir model and the production data. For two epochs a good comparison was obtained, for two others not. The stochastic error of the gravity values is small enough to detect the effect of gas extraction after a few years. But systematic errors present in the gravity data hamper the extraction of additional informational on the reservoir. Only with a very systematic network setup and well calibrated instruments gravimetry can contribute to the modeling of the gas extraction process.

1 Introduction

In the past 18 years, four gravity campaigns were carried out above the Groningen gas field. The campaigns were not set up for the prime purpose of reservoir modeling: the initial drive was the possible replacement of the costly leveling campaigns for the monitoring land subsidence by gravimetry. Although this did not work out, the (limited) gravimetry campaigns were continued. In 1996 the last campaign was carried out already with the idea in mind to find out what gravimetry could contribute to the modeling of the gas extracting process.

The Groningen gas fields are located in the northern part of the Netherlands. It is one of the world's bigger gas fields with an area of the main field of approximately 900 km² and an initial reserve estimated 2900 · 10⁹ m³. The mean depth of the layer is 2900 meters below sea level. The production started in 1963 and will be continued far into the next century. Apart from the main field, the *Groningen field*, many small subfields exist; mainly west and north-west to the Groningen field were the gas-carrying layer is rather fractured. See figure 1 for a map of the gas field.

In this paper we report on the gravimetry work and we would like to answer the following two questions:

- What is the gravity signal of the natural gas extraction?
- Can gravity surveys contribute to the reservoir modeling?

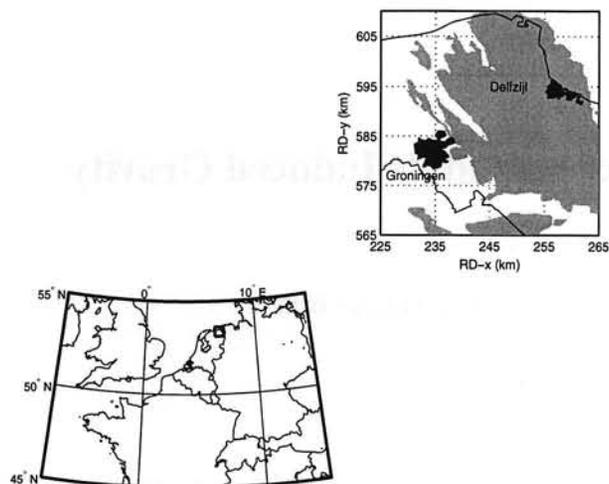


Fig. 1. The location of the Groningen gas fields.

First we compute the expected gravity signal from the production data. Then we see whether the gravity change from the reservoir data and the observed gravity change could give us hints about the correctness of the applied reservoir model.

In this paper we will first focus on the reservoir and production data that is available and how that relates to gravity. Then we will show some data and results of the gravimetry campaigns. In the following section we will confront the calculated gravity variations with the observed values and finally we will come to some conclusions and recommendations.

2 The Reservoir Data

In the Groningen field 26 production clusters are located. In these points the following information was available to us:

- depth of the gas-carrying layer
- thickness of the layer above gas-water contact
- proportion of natural gas (or Equivalent Hydro Carbon heights)
- pressure change for each epoch
- temperature

Furthermore, also for some inspection wells the depth of the layer was given and there is data about the composition of the natural gas. Each of the data listed above was interpolated to a denser point set by means of inverse distance interpolation. See figures 2 and 3. The thickness and therefore also the EHC values are rather smooth, whilst the depth is more irregular due to the fracturing of the reservoir layer. In the period 1978-1984 the pressure drop in the north is clearly higher than in the southern part of the reservoir. In that period more gas was produced from the northern area.

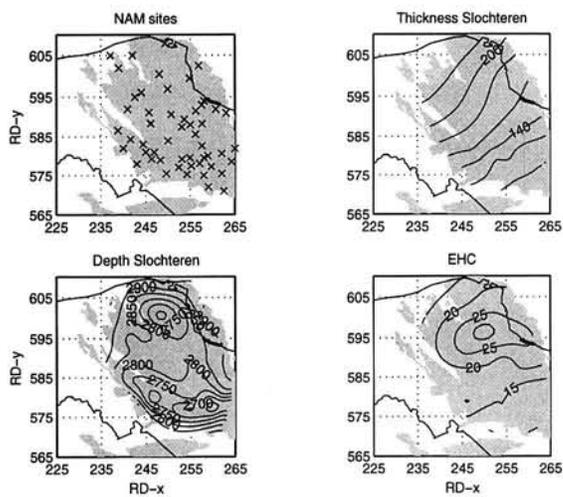


Fig. 2. Characteristics of the Groningen field.

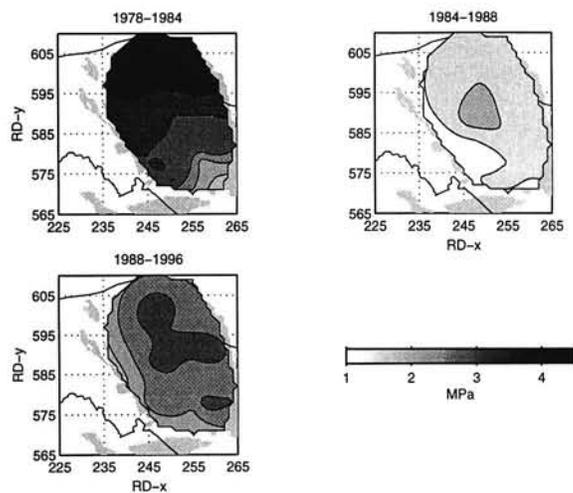


Fig. 3. Pressure changes for the three epochs.

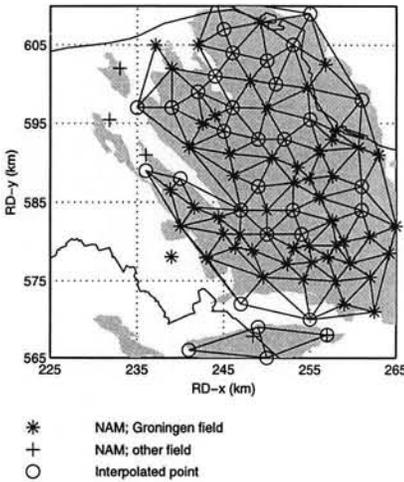


Fig. 4. The finite element model of the reservoir.

From this data the induced gravity reduction at the surface can be calculated. First density changes are derived from the production data as this can be directly related to gravity. For the gas density change we have

$$\Delta\rho = \frac{m\Delta n}{V_{gas}}$$

with m the molecular weight of the gas, n the number of mole gas and V_{gas} the volume of the gas. With the ideal gas law

$$PV_{gas} = zRTn$$

(P : pressure, z : correction factor, R : gas constant and T : temperature) this can be written as

$$\Delta\rho = \frac{m}{RT} \Delta \frac{P}{z}$$

By a volume integral over the reservoir the effect on the vertical component of the gravity vector at the surface is found:

$$\Delta g = G \iiint \frac{z}{\ell} \frac{\Delta\rho}{\ell^2} dV \approx G \iint \frac{z}{\ell} \frac{\Delta\rho}{\ell^2} EHC dA.$$

Obviously, the volume integral only covers the gas volume and not the entire layer thickness. The effective gas layer thickness is represented by the EHC values. As it is low with respect to the depth of the reservoir (30 vs. 2900 meters) the volume integral can be replaced by a surface integral over the reservoir layer. In theory the pressure should be corrected for reservoir compaction. The maximum surface subsidence amounts to 20 centimeters; with an average EHC of 20 meters this is an effect in the order of 1% and therefore negligible.

The latter integral was evaluated by numerical integration. A finite element model of the reservoir was constructed with the NAM wells as nodes plus an additional amount of interpolated support points; see figure 4. The attraction change for each epoch, triangle and gravity station combination was computed by a linear interpolation of the EHC values in the triangle and a Gauss quadrature of the surface integral (figure 5).

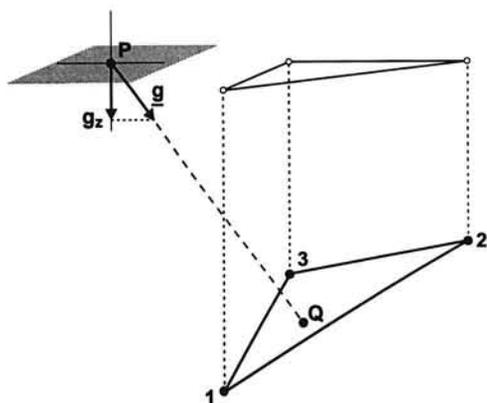


Fig. 5. Illustration of the integration over a triangle.

An independent check on the reservoir model was obtained by using the production data. The reservoir volume directly follows from the finite element model but can also be estimated from the combination of pressure and production data with the gas law:

$$\hat{V}_{gas} = \frac{z}{\Delta P} R \bar{T} \Delta n.$$

(The overbars indicate the average over the reservoir.) In this estimate it is assumed that the reservoir is rather homogeneous. The comparison yielded a difference of 12% which seems to be acceptable for the accuracies we strive for. (As such the total reservoir volume does not enter into the formulae at all, it serves only as an indication if we are on the right track with the reservoir model.)

The computed gravity changes for the period 1978-1996 are displayed in figure 6. The picture is rather smooth with a maximum decrease of $40 \mu\text{gal}$. Amplitudes which are well detectable with terrestrial gravimetry.

3 The Gravity Surveys

Four gravity campaigns were carried out in a 18 year time span: 1978, 1984, 1988 and 1996; an average interval of six years. The initial network consisted of 21 stations; 20 in the Groningen area and one outside as base point (see figure 7). Most of the stations are situated on NAM sites (NAM is the oil company exploiting the field). Some others are in railway stations or near churches. Mandatory for each point was the availability of leveling heights to be able to carry the free air reduction (see furtheron).

In 1984 and 1988 one new station was added to the network; mainly to get a better connection between the different epochs. In 1996 two more base stations were added. In each campaign 85 up to 106 measurements (i.e. gravity differences between stations) were made. For all the campaigns LaCoste Romberg G-gravimeters were used. Unfortunately for each campaign different instruments were available and, except for the last campaign, without an feedback system. This makes the observations prone to systematic errors which are hard to detect. According to the manufacturer, or see e.g. Becker (1984), Torge (1989) or Valliant (1991),

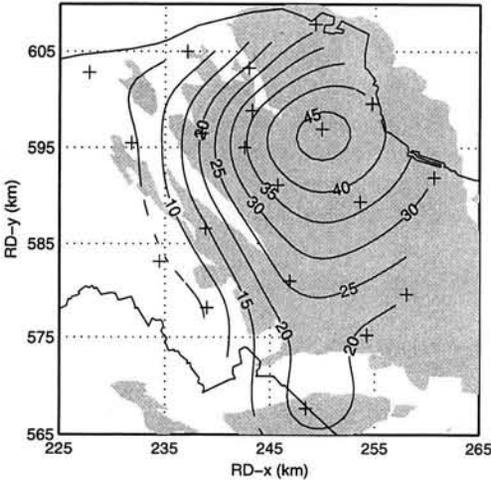


Fig. 6. The calculated gravity change in the epoch 1978 - 1996.

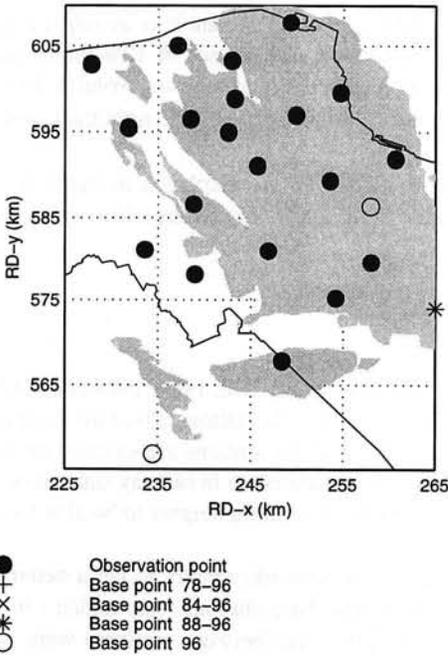


Fig. 7. The gravity network

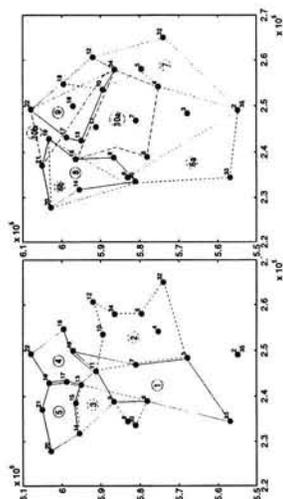


Fig. 8. The gravity network of the 1996 campaign. The network consists of 10 loops and 149 observations.

the G-type instruments can have periodic errors with an amplitude as high as $30\mu\text{gal}$. The random error for this kind of instruments has a standard deviation of $5 - 10\mu\text{gal}$ if all the necessary corrections have been applied.

For the 1996 network care was taken to get a robust and accurate network. Therefore a network planning was carried out to optimize the design. Essentially this is an error propagation using an a-priori instrument standard deviation and the calculation of the threshold values for the statistical testing. The following criteria were applied:

- maximum precision and reliability,
- inclusion of three points from the Dutch reference network,
- closed loops within the day for drift control,
- two observers/two instruments: each point visited by each possible combination,
- realistic travel time between the points,
- NAM site occupations concentrated in time as much as possible,
- maximum campaign duration of 12 days.

This led to the network shown in figure 8.

After the least-squares adjustment various statistical test were carried out to detect blunders in the readings and breaks in the instrument drift. Finally the network was connected to the base points to get absolute gravity values. Unfortunately it turned out that the newly added base points could not be used because they are close to the gas field and the information of the Dutch gravity network was too old. The formal, a-posteriori error standard deviations for each station and each campaign is shown in figure 9. In 1978 and 1984 the standard deviation is $6 - 7\mu\text{gal}$ and for the latest epochs $4\mu\text{gal}$. It has to be mentioned that especially for the epochs of 1978 and 1984 different data adjustment strategies led to significant different results for some points.

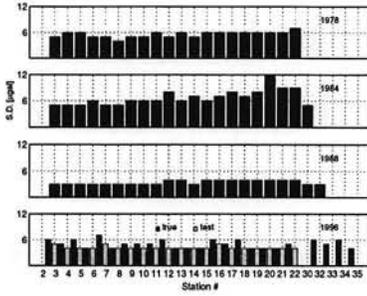


Fig. 9. The a-posteriori standard deviations of the gravity values for each station and epoch. The bars 'test' for 1996 refer to the values predicted with the network design.

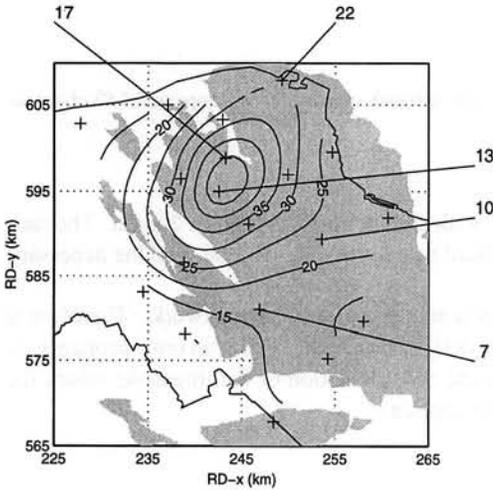


Fig. 10. The land subsidence for the epoch 1978 - 1996.

Also this indicates that the observational model is not completely correct (i.e. the data contains systematic errors)

The final correction that has to be applied is for land subsidence. Due to the compaction of the reservoir the surface subsides up to almost a centimeter per year. For the 1978 - 1996 period the subsidence is displayed in figure 10. The leads to a considerable effect in the surface gravity values ($3\mu\text{gal}/\text{cm}$). As we are interested in what is happening under the surface, all the absolute gravity values are reduced to the 1978 altitude by means of a free-air reduction. The final gravity changes for the epoch 1978 - 1996 are displayed in figure 11.

For a selected number of stations the values for all epochs are displayed in figure 12. The a-posteriori standard deviations for the gravity changes are about $\sqrt{6^2 + 4^2} \approx 8\mu\text{gal}$. Apart from systematic errors, these values are low enough to see the gravity change in one epoch (6 years) or e.g. the effect of ground water influx. (E.g., an area of 4 kilometer radius with two meter water at the depth of the reservoir yields about $15\mu\text{gal}$.) Stations 1, 2 and 4 are located in the center of the field, 5 is at the eastern border and 6 is outside the field to the north-west.

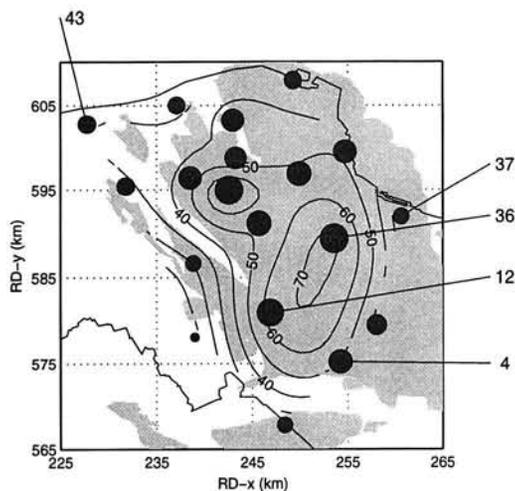


Fig. 11. The observed gravity changes for the period 1978 - 1996.

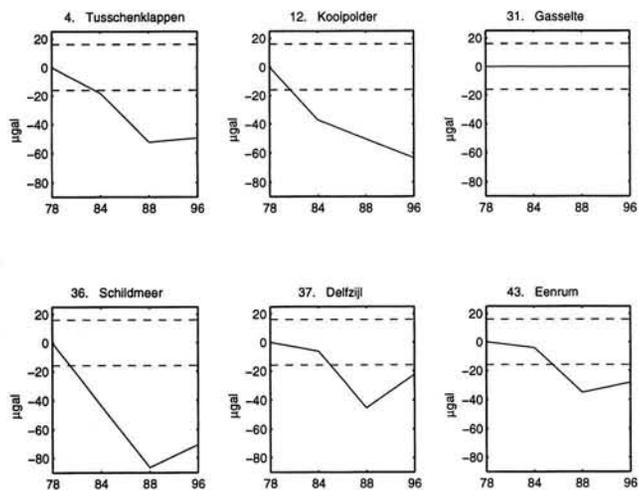


Fig. 12. The observed gravity changes for a selected number of stations. The numbers in front of the station names refer to the numbers in figure 11. The dashed line indicates the 2σ values.

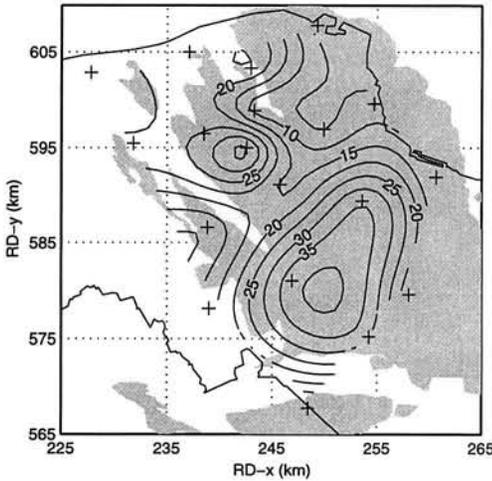


Fig. 13. Difference between observed and calculated gravity changes for the epoch 1978 - 1996.

A clear downward trend in gravity is visible; especially in the center of the field as could be expected. However, for 1988 a strange dip occurs. It seems that the gravity values of 1988 are too low. A careful re-examination of the 1988 campaign did not reveal any undetected errors and this effect remains unexplained up to now. The station of Gasselte (the basepoint) was only included for later reference. By definition at this station the gravity change is zero.

4 Confrontation of the Results

It could already be seen from the figures 6 and 11 that the results are quite different. In figure 3 the difference between the observed and calculated gravity changes is depicted. The a-priori standard deviation of the gravity observations was about $8\mu\text{gal}$. If we add a few μgal for the reservoir mismodelling this means that 2σ is about $20\mu\text{gal}$. About half of the number of points exceeds this threshold. All the observed gravity changes are too high with respect to the calculated change. As no probable geophysical scenario exist for this situation (if it was the other way round we could think of e.g. ground water influx) the deviations have to be attributed either to reservoir mismodelling or to undetected (systematic) errors in the gravity observations. For the four points in the north-west part of the area the former explanation is not improbable: in this area the reservoir is rather fractured and many small fields exist. This was not put accurately into the finite element model. For the point in Groningen city it is believed an height error exists for the 1978 and should be eliminated. The errors for the four points in the middle of the reservoir can only be attributed to gravimetry errors.

As was already mentioned in the introduction, the gravimeters used in the first three epochs were different and did not have a feedback system. Moreover, the network setup was different for each campaign. Although the gravimeters have been calibrated for a linear scale factor, no extended testing was performed to detect e.g. periodic errors. All this together makes the dataset prone to systematic errors which can be much larger than the stochastic measurement noise.

In figure 14 we displayed the difference between 1984 and 1996. These results show a good

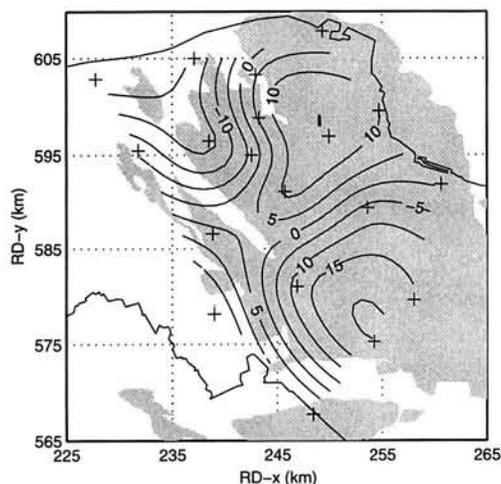


Fig. 14. Difference between observed and calculated gravity changes for the epoch 1984 - 1996.

agreement between the calculated and observed values. None of the points (except for point 6 in the NW part) exceeds the $20\mu\text{gal}$ threshold. This means that the gravimetry campaigns are of acceptable quality but also that the applied reservoir model is in agreement with these data and does not have to be improved.

In the next figure (15) the gravity data for six selected stations has been plotted. Here 1984 serves as reference. Clearly for some of the points displayed here the 1978 and 1988 do not give a good fit (i.e. the calculated gravity changes exceed the error bounds of the observations). The strange dip in 1988, as observed in the gravimetry data, does not occur in the computed gravity change. The point Gasselte, the base point, fortunately appears to be stable in time.

5 Conclusions and Recommendations

The conclusions of this research can be summarized as follows:

- A reasonable fit of observation data for epoch 1984 - 1996 has been obtained.
- Gravity data for 1978 and 1988 are questionable with respect to the gas data.
- In general the gravity data quality is too poor to yield additional information about the gas field. At best, a reasonable fit was obtained between observed and calculated data. This can be mainly attributed to undetected, systematic effects in the gravity data.

As the stochastic noise level of the instruments is very low, we still believe that gravimetry can provide a valuable contribution to the monitoring of gas fields. In order to do so we have the following recommendations:

- consistent measurement setup
- improved instrument calibration

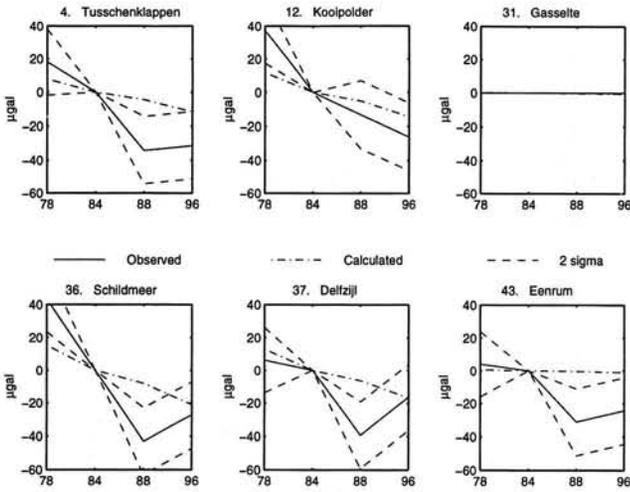


Fig. 15. The gravity changes for a number of selected stations.

- more epochs (minimal every 3 years or so) to be able to remove an epoch if something went wrong with the observations
- more stations needed for higher reliability of the detection of model errors
- more datum points desired for better control.

References

- Becker, M. (1984). Analyse von hochpräzisen Schweremessungen. Deutsche geodätische Kommission, Reihe C, Heft 294.
- Torge, W. (1989). *Gravimetry*. De Gruyter.
- Valliant, H. (1991). Gravity meter calibration at lacoste and romberg. *Geophysics*, 5, 705–711.

Error Propagation for Satellite Gradiometry

Martin van Gelderen

Abstract

Two approaches for the estimation of potential coefficients from satellite gradiometry are compared: least-squares and quadrature. Three errors are computed for each method: the propagated error, the bias and the aliasing error. Special attention is given to the effect of the regularization required by the presence of data gaps in the polar areas.

1 Introduction

For error prediction of a satellite gradiometer mission we often rely on (co)variance propagation, with degree variances or the full co-variance matrix, in the time or in the frequency domain. Most attention we gave to the propagated noise. With a model for the instrument noise, the variances of the potential coefficients can be simply estimated from the linear observation model. Two other main error sources, biases and aliasing, we usually omit because we believe they are small or they are more difficult to model. In this paper it is attempted to give some rough estimates of these two errors in particular in the view of the existence of polar gaps. The general idea is taken from the paper of Xu (1992), meanwhile some work in the same direction has been carried out by Bouman and Koop (1996).

2 The Model

Our calculations are carried out with the place domain approach. The general conclusions, however, will also be valid for the time domain approach because for dense data sets both methods will yield identical results. For the simplicity of the calculations and the clarity of the expressions the radial (zz) component was analyzed. For other components the result will differ but as we do not strive for accurate results but for getting a general impression this should not limit the validity of the conclusions.

We start with the linear observational model (everything in spherical, constant radius approximation)

$$z(\theta, \lambda) = \sum_{m,l} e^{im\lambda} P_{lm}(\cos \theta) (l+1)(l+2) \left(\frac{R}{r}\right)^{l+3} C_{lm}, \quad (1)$$

or in matrix notation

$$z = EP\Lambda c, \quad (2)$$

$$\underbrace{\begin{pmatrix} y_{m_1}(\theta_1) \\ \vdots \\ y_{m_1}(\theta_M) \\ y_{m_2}(\theta_1) \\ \vdots \\ y_{m_2}(\theta_M) \\ \vdots \end{pmatrix}}_y = \underbrace{\begin{pmatrix} P_{0,m_1}(\theta_1) & \cdots & P_{L,m_1}(\theta_1) & | & 0 & \cdots & 0 & | & 0 \\ \vdots & & \vdots & & \vdots & & \vdots & & \vdots \\ P_{0,m_1}(\theta_M) & \cdots & P_{L,m_1}(\theta_M) & | & 0 & \cdots & 0 & | & 0 \\ \hline 0 & \cdots & 0 & | & P_{0,m_2}(\theta_1) & \cdots & P_{L,m_2}(\theta_1) & | & 0 \\ \vdots & & \vdots & & \vdots & & \vdots & & \vdots \\ 0 & \cdots & 0 & | & P_{0,m_2}(\theta_M) & \cdots & P_{L,m_2}(\theta_M) & | & 0 \\ \hline 0 & \cdots & 0 & | & 0 & \cdots & 0 & | & 0 \end{pmatrix}}_P \underbrace{\begin{pmatrix} C_{0,m_1} \\ \vdots \\ C_{L,m_1} \\ C_{0,m_2} \\ \vdots \\ C_{L,m_2} \\ \vdots \end{pmatrix}}_x$$

$$\underbrace{\begin{pmatrix} z(\theta_1, \lambda_1) \\ \vdots \\ z(\theta_1, \lambda_N) \\ z(\theta_2, \lambda_1) \\ \vdots \\ z(\theta_2, \lambda_N) \\ \vdots \end{pmatrix}}_z = \underbrace{\begin{pmatrix} e^{im_1\lambda_1} & 0 & \cdots & | & e^{im_2\lambda_1} & 0 & \cdots & | & \vdots \\ \vdots & \vdots & & & \vdots & \vdots & & & \vdots \\ e^{im_1\lambda_N} & 0 & \cdots & | & e^{im_2\lambda_N} & 0 & \cdots & | & \vdots \\ \hline 0 & e^{im_1\lambda_1} & 0 & \cdots & 0 & e^{im_2\lambda_1} & 0 & \cdots & \vdots \\ \vdots & \vdots & \vdots & & \vdots & \vdots & \vdots & & \vdots \\ 0 & e^{im_1\lambda_N} & 0 & \cdots & 0 & e^{im_2\lambda_N} & 0 & \cdots & \vdots \\ \hline \cdots & & & & \cdots & & & & \vdots \end{pmatrix}}_E \underbrace{\begin{pmatrix} y_{m_1}(\theta_1) \\ \vdots \\ y_{m_1}(\theta_M) \\ y_{m_2}(\theta_1) \\ \vdots \\ y_{m_2}(\theta_M) \\ \vdots \end{pmatrix}}_y$$

Fig. 1. The structure of the matrices (Λ not shown here: it simply is diagonal)

where c contains the potential coefficients, Λ the eigenvalues related to the type of observations, P the Legendre functions, E the exponential base functions and z the observations. The structure of the matrices is illustrated in figure 1. Multiplying both sides by $\frac{1}{N}E^*$ (conjugate transpose, N number of observations on a parallel) yields the new observation equations

$$y \equiv \frac{1}{N}E^*z = P\Lambda c. \quad (3)$$

Due to the block-diagonal structure of P and its orthogonality properties the solution can be computed for each order and parity of degree individually. We'll concentrate on this observation equation. Although it does not give the complete picture, the main disturbing effect is the presence of polar gaps which enter only in this step.

Two kind of solutions will be considered: least-squares and quadrature of the corresponding inverse formula. Generally the solution can be written as

$$\hat{c} = B^T y.$$

For the *least-squares solution* we have

$$B^T = (P^T Q_{yy}^{-1} P + R)^{-1} P^T Q_{yy}^{-1},$$

with a (diagonal) regularization matrix R . This matrix is required for the lower orders: due to the polar gaps the normal matrix gets (weakly) singular. Kaula's rule was used here for R .

For the radial component it is easy to find the analytical inverse of (3) in the limit case of an infinite number of observations and the whole earth covered homogeneously:

$$C_{lm} = \frac{1}{(4 - 2\delta_{m,0})} \frac{1}{(l+1)(l+2)} \left(\frac{r}{R}\right)^{l+3} \int_{\sigma} y(\theta) P_{lm}(\theta) \sin \theta d\theta. \quad (4)$$

For the case with real data, the integral is discretized in the area where data is available. In matrix notation:

$$B^T = \Lambda^{-1} P^T S$$

(S is a diagonal matrix with $\sin \theta$ and the necessary scale factors).

3 The three errors

The complete picture of the error (disregarding the approximations committed by using the model (1) as starting point) consists of three parts: the propagated error, the bias and the aliasing error. The *propagated error* is

$$B^T Q_{yy} B, \quad (5)$$

where Q_{yy} is the covariance matrix of the observation noise. The *bias* is defined as

$$(I - B^T P \Lambda), c$$

but this requires the (true) potential coefficients. To avoid this, and to make the calculation easier, we compute the signal power of the bias (*bias-variance*):

$$(I - B^T P \Lambda) K (I - B^T P \Lambda)^T. \quad (6)$$

Here K is a diagonal matrix with degree/order variances of the potential, e.g. from OSU91a or Kaula. The third error is introduced when discretizing the observation equations: meanwhile the number of unknowns (potential coefficients) had to be reduced in order to avoid an under-determined system of equations. This introduces a model error (*aliasing*), whose effect can be numerically estimated by extending the model (2) as

$$y = \begin{pmatrix} P & P' \end{pmatrix} \begin{pmatrix} \Lambda & 0 \\ 0 & \Lambda' \end{pmatrix} \begin{pmatrix} c \\ c' \end{pmatrix},$$

where P , Λ and c are defined as before and their primed versions complete the model up to degree infinity. Infinity is here approximated by $L = 360$ (the maximum degree and order of the coefficients estimated is 180). Now the total non-stochastic error is computed as

$$c - \hat{c} = c - B^T y = c - B^T P \Lambda c - B^T P' \Lambda' c'.$$

The first two terms together we already defined as the bias, the last term is the *aliasing* error. Also this error we model with signal variances as:

$$(B^T P' \Lambda' c') K (B^T P' \Lambda' c')^T.$$

Again we underline that this is not the complete aliasing error: implicitly we assumed that in the east-west direction no aliasing occurs.

4 Results

For the computations a noise PSD of the gradiometry data $10^{-3}\text{E}/\sqrt{\text{Hz}}$ and an orbital inclination of 97.5° were used. For the noise we took the model of Jekeli-Rapp Jekeli and Rapp (1980):

$$\sigma_{lm}^2 = \frac{\sigma_{obs}^2}{N \cdot M}.$$

The parameters were not tuned very accurately. It was assumed that there is (on the average) one observation in a half-by-half degree block (N and M denote the number of observations per parallel and meridian, respectively), which compares to a sampling time of approximately 10 seconds. The height of the orbit is 350 kilometers.

The results of the analysis are shown in figure 2. For some orders the results are also presented as line graphs (figure 3). It can be clearly seen that least-squares has a superior performance for the propagated error, as expected. There is a notable bias in order zero with quadrature whereas for LS this is more spread over the spectrum. The main point of concern, however, is the considerable aliasing error in the least-squares estimate. (We should still keep in mind that the sampling effect in longitudinal direction was not accounted for here.)

To estimate the effect of the selection of the maximum degree of the linear observation model, we take $L = 90$ (figure 4). For quadrature only the aliasing error changes, but with least-squares the minimum degree affected by aliasing shifts with the maximum degree to the left. This leads to the conclusion that increasing the maximum degree in the linear model might diminish the aliasing error. This would mean that, if possible with the data sampling, a maximum degree of, e.g. 240 is used instead of 180 in the LS procedure, the aliasing error can be reduced, although the signal-to-noise ratio of the coefficients above 180 might be smaller than one (henceforth do not contain useful information) and are thrown away afterwards. Another observation is the decrease of the bias for least-squares. This might be explained by the more strict band-limitation. This reduces the singularity of the normal system and has the effect that the implicit extrapolation of the measurements into the polar gaps gets more weight in the least-squares adjustment (see also discussion in the next section).

The largest errors in the spectrum are always in a wedge-shaped area of the spectrum. This can be explained by the power-loss of the Legendre functions due to the polar gaps Gelderen and Koop (1997), by calculating the condition number of the normal matrix or by simply inspecting the Legendre functions. As can be seen from figure 5, the support of the Legendre functions in the polar areas decreases for an increasing orders and decreasing degrees.

Another typical pattern we see is that within this wedge, the error increases with increasing l up to degree 150 and decreases a little going towards 180. This behavior can be explained by doing an eigenvalue analysis of the normal matrix. The eigenvalue decomposition of the normal matrix is:

$$N = P^T \Lambda Q_{yy}^{-1} \Lambda P \equiv SDS^{-1}.$$

Transforming the stabilized, inverse normal matrix to this system of eigenvectors yields:

$$\Rightarrow S^{-1}(N + R)^{-1}S = (SNS^{-1} + SRS^{-1})^{-1} = (D + SRS^{-1})^{-1}. \quad (7)$$

Clearly, for small eigenvalues of N , the regularization matrix R gets all the weight in the inversion. These components will be badly estimated because they are almost completely determined

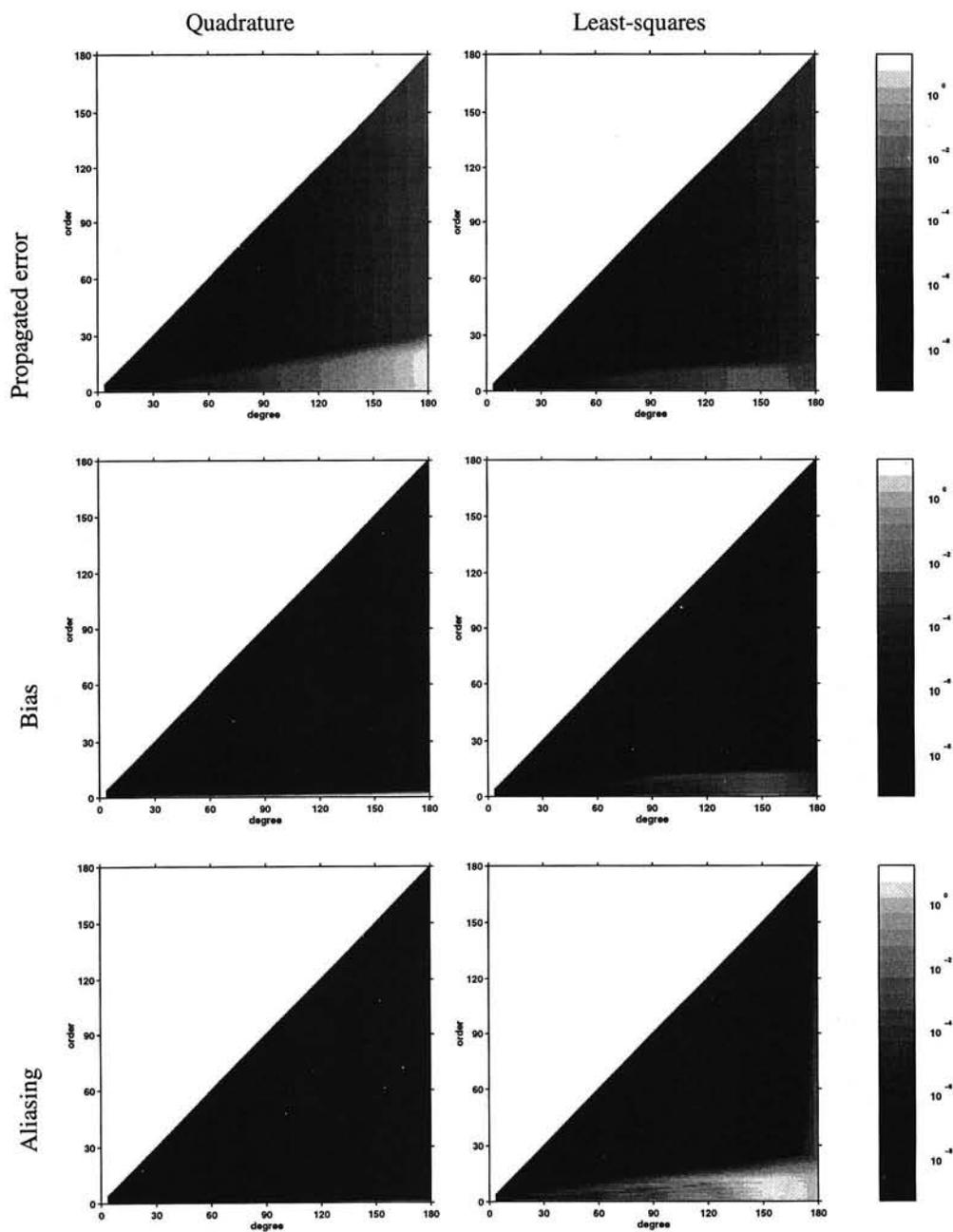


Fig. 2. The results for all potential coefficients

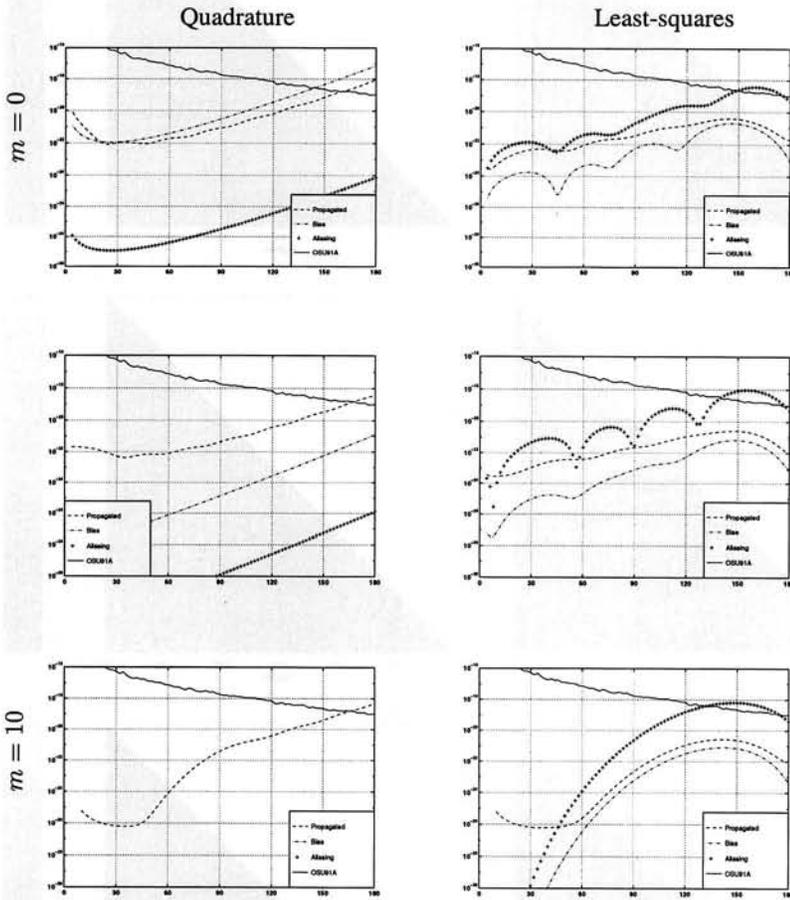


Fig. 3. Results for some orders in particular

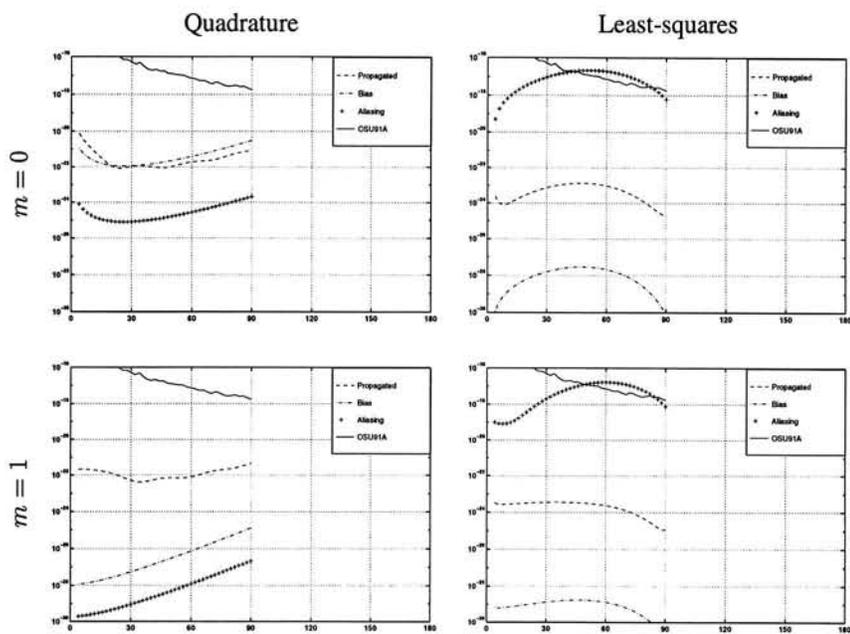


Fig. 4. Results for the orders zero and one for $L = 90$

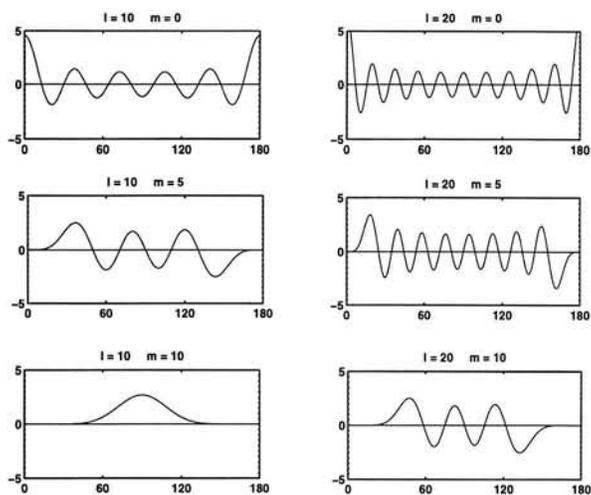


Fig. 5. Some Legendre functions (horizontal axes: co-latitude in degrees)

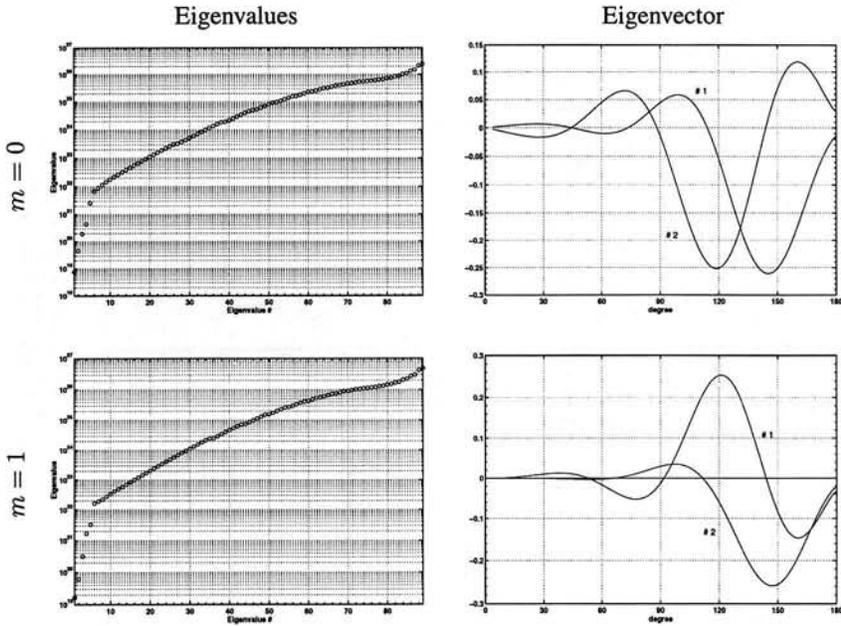


Fig. 6. The eigenvalues of the normal matrix and the eigenvectors of the smallest eigenvalues

by the regularization. The components of the related eigenvectors indicate the linear combinations of coefficients which are estimated poorly. For many cases, these eigenvectors have a similar character: increasing support towards \pm degree 150, then decreasing. Two examples are shown in figure 6.

5 Contribution of Regularization

Regularization is necessary to obtain a solution when the data is not sufficient to solve for the parameters directly. We would like to see how much is the effect of the regularization on our final estimate. The contribution of the regularization to the least-squares solution can be estimated as follows. We start with the standard partitioned model of observation equations:

$$\begin{pmatrix} z_1 \\ z_2 \end{pmatrix} = \begin{pmatrix} A_1 \\ A_2 \end{pmatrix} c, \quad W = \begin{pmatrix} W_1 & 0 \\ 0 & W_2 \end{pmatrix}.$$

If c can be estimated from z_1 and z_2 individually (\hat{c}_1 and \hat{c}_2 respectively), we can write

$$\hat{c} = \underbrace{(A_1^T W_1 A_1 + A_2^T W_2 A_2)^{-1}}_{N^{-1}} \underbrace{(A_1^T W_1 A_1)}_{N_1} \hat{c}_1 + \underbrace{A_2^T W_2 A_2}_{N_2} \hat{c}_2 \quad (8)$$

The relative contribution of the observations z_1 to the final estimate is $N^{-1}N_1$, like-wise for z_2 . The standard method of regularization is obtained by setting

$$\begin{aligned} A_2 &= I, & z_2 &= 0, \\ W_2^{-1} &= \text{degree/order variance model}, \end{aligned}$$

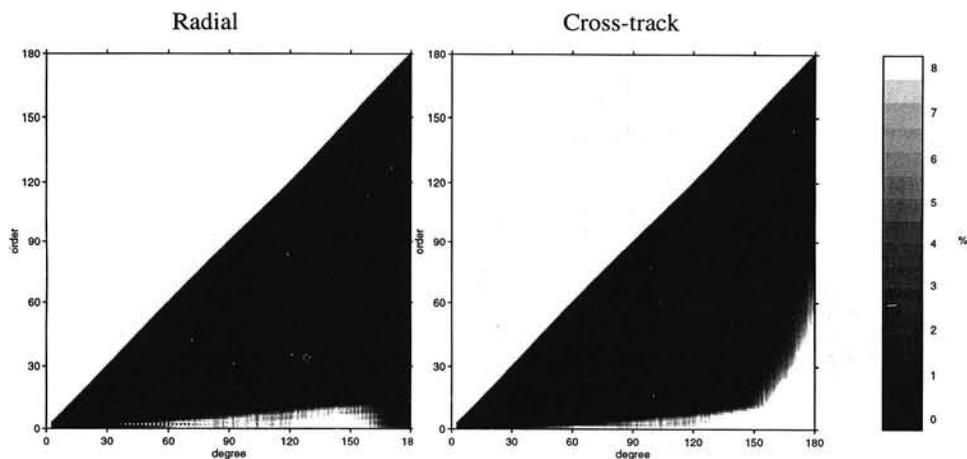


Fig. 7. Relative contribution of regularization

and obtain

$$\hat{c} = N^{-1}(N_1\hat{c}_1 + W_2\hat{c}_2).$$

The relative contribution of regularization is then $N^{-1}W_2$. For the radial and cross-track component the results are shown in figure 7. The contribution of regularization is mainly concentrated in the low orders. The cross-track component is weak in the low-order/high-degree part of the spectrum, which is also reflected in the relative contribution of the regularization.

As alternative for the standard regularization we can use zero observations only for the polar areas ('local regularization'). Then we take

$$W_2 = I - W_1$$

(we assume here all observations have unit weight). If the radial component has been observed the design matrix can be written as

$$A = Y\Lambda,$$

($Y = EP$) and the relative contribution of the estimation \hat{c}_2 (the estimation from the fake observations) becomes

$$N^{-1}N_2 = \Lambda^{-1}Y^T P_2 Y \Lambda.$$

The diagonal elements of this expression are

$$\int_{\text{Polar gaps}} Y_{lm}^2 d\sigma$$

and do not depend on the type of data observed! It is equal to the power-loss rule mentioned earlier; see figure 8. (Actually \hat{c}_1 and \hat{c}_2 can't be computed at all because both are given on a part of the earth's surface only, but here we assume both z_1 and z_2 are globally available only with a very low weight for the poles or the remaining domain respectively).

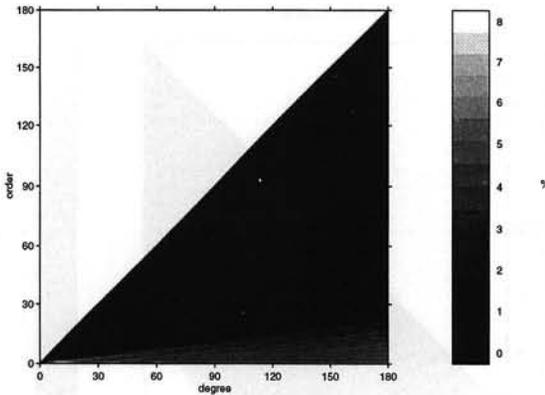


Fig. 8. Relative contribution of local regularization

As we are mainly interested in geoid heights etc., the relative contributions are propagated to the space domain. Suppose we generate a grid of, e.g. geoid heights from the estimated coefficients. If the grid is dense enough we can compute back the coefficients without loss of accuracy. In matrix form:

$$f = Gc \equiv Y\Lambda_f c \quad c = \Lambda_f^{-1} Y^{-1} f.$$

(f contains the grid values, G is the transfer matrix). f can be estimated from \hat{c} , \hat{c}_1 or \hat{c}_2 which we use to compute the relative contribution of regularization

$$\hat{f} = G\hat{x}, \quad \hat{f}_1 = G\hat{c}_1, \quad \hat{f}_2 = G\hat{c}_2$$

$$\Rightarrow \hat{f} = Y\Lambda_f N^{-1} N_1 \Lambda_f^{-1} Y^T \hat{f}_1 + \underline{Y\Lambda_f N^{-1} N_2 \Lambda_f^{-1} Y^T \hat{f}_2}.$$

For standard regularization the relative contribution of regularization (the underlined part) is

$$Y\Lambda_f N^{-1} P_2 \Lambda_f^{-1} Y^T.$$

For the local method it is

$$Y \Lambda_f \Lambda^{-1} Y^T P_2 Y \Lambda \Lambda_f^{-1} Y^T$$

The results are shown in figure 9. Contrary to what was expected, the local regularization performs worse with respect to the standard regularization: it has more influence on the computed geoid etc. It can be improved a little by smoothing the transition from true values to zero values at the polar gap boundary by a linear filter (figure 10), but that does not change the results fundamentally. The reason for the weaker performance is that though the regularization is very local, it is also hard: the fake data gets all the weight in the polar zones. With the standard method it is much softer: we put in zero potential coefficients but with a high variance. This makes that this fake information is also used very selective due to the least-squares algorithm. What it really makes performing so well is the band-limitation we apply in the linear model. Together with the high amount of data available the least-squares procedure extrapolates to the polar areas and the relative contribution of regularization never reaches 100% in the polar gaps. Whether this good performance is real, after all the band-limitation is also artificial, should be investigated further.

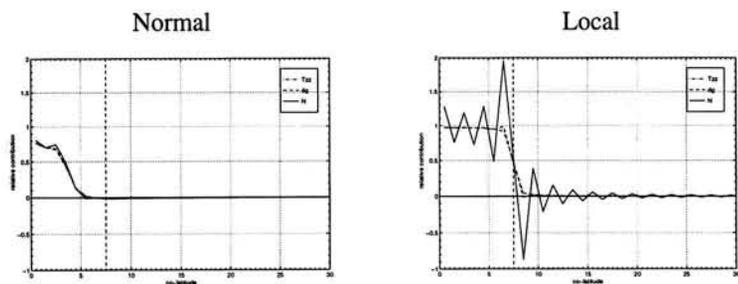


Fig. 9. Relative contribution of regularization propagated into space domain

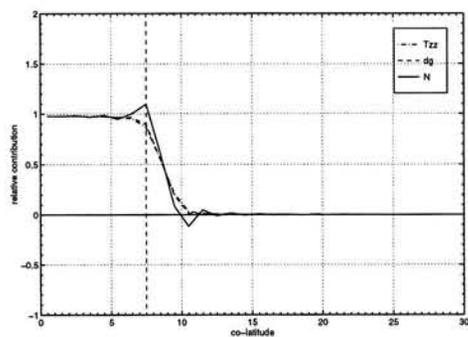


Fig. 10. Relative contribution of local regularization with filter propagated into space domain

6 Conclusions & Recommendations

- Least-squares gives a smaller over-all error with respect to quadrature, mainly due to the smaller propagated error.
- The bias introduced by regularization is not very large but that might be related to the band-limitation.
- The aliasing error can exceed the signal below degree 180.
- The local regularization does not work satisfactory.
- The effect of the true bias and aliasing error should be investigated by means of simulation studies. The maximum degree used in the observation equations is of particular importance due to the non-orthogonality of the base-functions in the observed area.

References

- Bouman, J. and Koop, R. (1996). Regularization in gradiometric analysis. pres. at EGS The Hague.
- Gelderen, M. V. and Koop, R. (1997). The use of degree variances for gradiometric analysis. *J. of Geod.*, **71**, 337-343.
- Jekeli, C. and Rapp, R. (1980). Accuracy of the determination of mean anomalies and mean geoid undulations from a satellite gravity field mapping mission. report no. 307, Ohio State University, Dept. of Geodetic Science.
- Xu, P. (1992). The value of minimum norm estimation of geopotential fields. *Geoph. J. Int.*, **111**, 170-178.

Quality Differences between Tikhonov Regularization and Generalized Biased Estimation in Gradiometric Analysis

Johannes Bouman and Radboud Koop

Abstract

The determination of the earth's gravity field from satellite gravity gradiometry is an inverse or ill-posed problem requiring regularization. In geodesy the regularization is usually done by adding a priori information, Kaula's rule for example. This can be interpreted as constraining the signal and is called Tikhonov regularization. Another method is to add an arbitrary positive definite matrix to the system of normal equations. This matrix is chosen such that the total deviation of the estimated function from the true function is minimal which is called biased estimation. The aim here is to compare the regularization by constraining the signal to biased estimation. This is done for several gradiometric mission scenarios. The comparison is based on the mean square error of the solutions, which is the sum of the propagated error and the, often neglected, regularization error. The total error of the solutions is computed and the errors are propagated to geoid heights as well. The main conclusions are that the qualification 'best regularization method' depends on the cause of the instability, that the regularization error is not negligible and that additional measurements or other solution methods are needed.

1 Introduction

An accurate and high resolution knowledge of the earth's gravity field is needed in several earth oriented sciences. In geodesy, for example, the gravity field is needed for levelling with GPS, in oceanography it is important for studying large scale ocean circulation and last but not least in geophysics a better knowledge of the earth's gravity field yields better boundary conditions in the study of the earth's interior.

The determination of the earth's gravity field is very convenient using satellite methods since a satellite orbiting the earth samples practically the whole globe within a relative short time span. Two very promising satellite techniques for global gravity field determination are satellite-to-satellite tracking and satellite gravity gradiometry. Here only gradiometry is considered since with this technique one can in principle determine all frequencies up to high degree and order, with the for gradiometry specific numerical instability caused by the polar gaps.

A disadvantage of both techniques is the ill-posedness of the problem, i.e. the solution, in particular gravity potential coefficients, derived from the measurements is unstable. One reason

is the downward continuation of the observations to the earth's surface. A stable solution can only be obtained by regularizing the solution. This is well known and often Kaula's rule is used, which can be interpreted as a constraint on the signal and this is called Tikhonov regularization (TR), Tikhonov and Arsenin (1977). An alternative method is to add an arbitrary positive definite matrix to the system of normal equations such that the difference between the estimated and true function is minimal in the sense of the mean square error. This is called generalized biased estimation (GBE) or ridge regression, Vinod and Ullah (1981); Xu (1992a, 1992b).

The purpose of this paper is to compare the two regularization methods for eight gradiometric missions. The comparison is based on the mean square error which is the sum of the propagated measurement error and regularization error or bias. The latter is often neglected, but is inherent to the regularization, Louis (1989); Xu (1992b).

The description of the eight example gradiometric missions and the observation model in Section 2 is followed by a summary of the methods of regularization and the related errors in Section 3. Section 4 lists the results and Section 5 presents the conclusions.

2 Model and mission description

2.1 Observation model

The unknowns to be solved for are the normalized harmonic coefficients \bar{C}_{lm} , \bar{S}_{lm} of a spherical harmonic expansion of the gravitational potential:

$$V = \frac{GM}{R} \sum_{l=0}^L \left(\frac{R}{r}\right)^{l+1} \sum_{m=-l}^l \bar{Y}_{lm}(\theta, \lambda) \quad (1)$$

with the abbreviation

$$\bar{Y}_{lm}(\theta, \lambda) = \begin{cases} \bar{C}_{lm} \cos m\lambda \bar{P}_{lm}(\cos \theta), & m \geq 0 \\ \bar{S}_{l|m|} \sin |m|\lambda \bar{P}_{l|m|}(\cos \theta), & m < 0 \end{cases} \quad (2)$$

where GM is the gravitational constant times mass of the earth, R the radius of a reference sphere enclosing all masses, l, m degree and order, $\bar{P}_{lm}(\cos \theta)$ the fully normalized Legendre functions and r, θ, λ the geocentric polar coordinates. For the maximum degree and order we take $L = 180$, corresponding to a spatial resolution of approximately 1° , which is a typical resolution to be achieved from a gradiometry mission.

The observations we consider are gravity gradients or the second order derivatives of the gravitational potential. The measurements could for example be the change in range between two falling proof masses around the earth. A local satellite coordinate system is x, y, z with x along-track, y cross-track and z radial. Observing the range changes in these three directions yields the observables V_{xx} , V_{yy} and V_{zz} . By a proper coordinate transformation these values can be related to (1), see e.g. Koop (1993).

The unknowns and observations are connected by the linear model

$$E\{g\} = Af, \quad D\{g\} = P^{-1} \quad (3)$$

with g the observations, f the unknowns, A the design matrix and P^{-1} the error covariance matrix of the observations. Let's assume that the orbit of the satellite is circular, that there are no data gaps and that after a number of revolutions the ground-track of the satellite repeats

Table 1. Missions considered

observation	inclination (degrees)				
	90	92.5	95	97.5	100
V_{yy}	+	+	+	+	+
V_{zz}	+			+	
V_{xx}, V_{yy}, V_{zz}				+	

itself exactly. Then one can consider the observations V_{xx} etc. as a time series along the orbit. Due to the assumptions one can compute the Fourier coefficients of these observations, the lumped coefficients, e.g. Koop (1993); Schrama (1990). These lumped coefficients are linear combinations of the unknown potential coefficients $\bar{C}_{lm}, \bar{S}_{lm}$. The above approach is the time-wise in the frequency domain method, with the advantage that for example colored noise can easily be accounted for, Rummel *et al.* (1993), see also Section 2.3.

2.2 Mission variables

Eight satellite gradiometric missions are considered as listed in Table 1.

V_{yy} was chosen since it was the observable for the proposed STEP mission, Blaser *et al.* (1996). Another ESA proposed mission, current under investigation, is GOCE (Gravity Field and Ocean Circulation Explorer), ESA (1996). This mission will measure the three diagonal components of the gravity tensor V_{xx}, V_{yy} and V_{zz} , for brevity denoted by V_d (diagonal). The observation V_d is used as reference. Several inclinations of the satellite orbit are considered.

2.3 Mission constants

For all missions we have chosen a satellite height of 300 km and a mission duration of six months. For the measurement error we took a colored noise PSD with a flat spectrum for $\beta_{km} > 10$ cpr at the level of $10^{-3} E/\sqrt{Hz}$ (cpr = cycles per revolution, $E = \text{Eötvös unit} = 10^{-9}/s^2$) and a $1/\omega$ characteristic for $2 \leq \beta_{km} \leq 10$ cpr. (ω stands for frequency here.) The β_{km} describe the spectrum along the orbit. The noise characteristic basically implies that the minimum degree that can be determined is 2, and that all spherical harmonic degrees above 2 are affected by the colored noise error spectrum.

3 Stabilization methods

3.1 Introduction

The standard technique for parameter determination is least-squares (l.s.). The sum of the squared errors has to be minimized:

$$\min_f \|g - Af\|_P^2 \quad (4)$$

which leads to the l.s. estimate \hat{f} of f

$$\hat{f} = N^{-1}A^T P g = A^+ g \quad (5)$$

where $N = A^T P A$ and P is the weight matrix of the observations g . However, the inverse of the normal matrix is unstable, reflecting an ill-posed problem. This has three reasons:

- *Downward continuation.* The higher degrees attenuate by a factor of $(R/r)^{l+1}$ at satellite height so the observation noise is amplified due to the downward continuation.
- *Polar gaps.* Every inclination not equal to 90° results in two polar gaps. Hence, a global gravity field is estimated from regional measurements.
- *Type of observation.* Every kind of observation related to the gravity potential (like gravity, position or gravity gradients) will have, in the frequency domain, a different sensitivity for different frequencies. For instance, for V_{zz} the sensitivity decreases with increasing l , whereas it is constant for all orders m per degree. Or V_{yy} which has an increasing sensitivity for increasing order m . Sometimes a particular observation is not sensitive to a certain gravity field parameter at all, like V_{yz} and V_{xy} in a polar orbit which are not sensitive to the zonal harmonics or V_{yy} from which the zonal harmonics can only poorly be determined, in particular at the equator.

Several methods exist to compute a stable solution. The first method discussed here is Tikhonov regularization, Tikhonov and Arsenin (1977), the second is ridge regression or biased estimation, Vinod and Ullah (1981); Xu (1992a, 1992b).

Using Tikhonov regularization (TR) implies constraining the signal f or some higher derivatives of f , whereas (generalized) biased estimation (GBE) just adds an arbitrary positive definite matrix to N to stabilize the solution. After introducing both methods in more detail, we compare the total Mean Square Error (MSE) in both cases.

3.2 Tikhonov regularization

Method. The idea is to constrain the total power of the signal or of derivatives of the signal. The minimization problem

$$\min_f \|g - Af\|^2 + \alpha \|Lf\|^2 \quad (6)$$

has to be solved instead of (4), where L is some differential operator and α is the compromise between minimizing the observation error and the constraint, Louis (1989); Engl *et al.* (1996).

Here we consider only the signal constraint with weighted norm

$$\min_f \|g - Af\|_P^2 + \alpha \|f\|_K^2 \quad (7)$$

so that

$$f_t = (N + \alpha K)^{-1} A^T P g = A_\alpha^+ g \quad (8)$$

where K is for example a diagonal matrix with elements $10^{10} l^4$ which is the inverse of the well known degree-order Kaula rule. If one takes for f a geopotential model, say OSU91A, then it approximately holds

$$f_t \approx \left(N + \alpha \left(\text{diag}(f f^T) \right)^{-1} \right)^{-1} A^T P g \quad (9)$$

which is needed for comparison with biased estimation.

Errors. Suppose one has exact observations g_e with corresponding solution $f = A^+g_e$. The difference between f and the regularized solution f_t is

$$f_t - f = A_\alpha^+(g - g_e) + (A_\alpha^+ - A^+)g_e. \quad (10)$$

The first term on the right hand side is called the data error, the second the regularization error, Louis (1989). The latter can be written as

$$\begin{aligned} \Delta f &= ((N + \alpha K)^{-1}A^TP - N^{-1}A^TP)g_e \\ &= ((N + \alpha K)^{-1} - N^{-1})Nf \\ &= (N + \alpha K)^{-1}(N + \alpha K - \alpha K)f - If \\ &= -(N + \alpha K)^{-1}\alpha Kf \end{aligned} \quad (11)$$

which represents the bias, cf. next section.

The total error or Mean Square Error Matrix (MSEM) consists of the propagated error, from (8)

$$Q_f = (N + \alpha K)^{-1}N(N + \alpha K)^{-1} \quad (12)$$

and the bias

$$MSEM = Q_f + \Delta f \Delta f^T. \quad (13)$$

The trace of the MSEM is called the Mean Square Error (MSE) which is used here as a measure of the quality of the solution.

3.3 Generalized biased estimation

Preliminaries. Biased estimation is also called ridge regression, e.g. Vinod and Ullah (1981). Let the Choleski decomposition of P be $P = WW^T$, and define the transformations $g_w = W^Tg$ and $A_w = W^TA$. Then the observation model is $E\{g_w\} = A_w f$ with least-squares solution

$$\begin{aligned} \hat{f} &= (A_w^T A_w)^{-1} A_w^T g_w \\ &= (A^T W W^T A)^{-1} A^T W W^T g \\ &= (A^T P A)^{-1} A^T P g = N^{-1} A^T P g. \end{aligned} \quad (14)$$

Let the singular value decomposition of A_w be

$$A_w = U \Sigma V^T \quad (15)$$

where U, V are orthogonal matrices, i.e. $U^{-1} = U^T$ and Σ is a diagonal matrix with singular values σ_i in descending order: $\lim_{i \rightarrow \infty} \sigma_i = 0$. The least-squares solution (14) can now be written as

$$\begin{aligned} \hat{f} &= (V \Sigma^2 V^T)^{-1} V \Sigma U^T g_w \\ &= V (\Sigma^2)^{-1} \Sigma U^T g_w \end{aligned} \quad (16)$$

since $A_w^T A_w = V \Sigma^2 V^T$ and $V^T = V^{-1}$. Now it is clear that the instability is caused by the small singular values σ_i , their inverse becomes large.

Method. The idea of generalized biased estimation is to add an arbitrary positive definite diagonal matrix to the matrix of singular values to stabilize the solution:

$$f_b = V(\Sigma^2 + D)^{-1}\Sigma U^T g_w. \quad (17)$$

Hence the smallest singular values become larger, the inverse can be computed.

Errors. The propagated error is, Xu and Rummel (1994):

$$Q_f = V(\Sigma^2 + D)^{-1}\Sigma^2(\Sigma^2 + D)^{-1}V^T. \quad (18)$$

The bias is, Vinod and Ullah (1981); Xu (1992b):

$$E\{f_b - f\} = \Delta f = -(N + M)Mf \quad (19)$$

where $M = VDV^T$ or

$$\Delta f = -V(\Sigma^2 + D)^{-1}DV^Tf. \quad (20)$$

Again the total error is the mean square error matrix, equation (13). The Mean Square Error is, Xu and Rummel (1994):

$$MSE = \sum_{i=1}^n \frac{\sigma_i^2 + d_i^2 \langle f, v_i \rangle^2}{(\sigma_i^2 + d_i)^2}. \quad (21)$$

The set of d_i with minimum MSE is obtained by differentiating the MSE with respect to d_i , see for example Xu and Rummel (1994):

$$\frac{\partial MSE}{\partial d_i} = \frac{2\sigma_i^2(d_i \langle f, v_i \rangle^2 - 1)}{(\sigma_i^2 + d_i)^3}. \quad (22)$$

The minimum is thus obtained for $d_i = \langle f, v_i \rangle^{-2}$.

Comparison with Tikhonov regularization. Note that

$$d_i^{-1} = \langle f, v_i \rangle^2 \quad (23)$$

where v_i is the i -th column of V . Hence we may write

$$D^{-1} = \text{diag}(V^T f f^T V) \quad (24)$$

which gives

$$f_b = \left(N + \left(V \text{diag}(V^T f f^T V) V^T \right)^{-1} \right)^{-1} A^T P g. \quad (25)$$

Comparing TR and GBE, equations (9) and (25), one sees some similarity in the regularization matrix. The difference is that the signal $f f^T$ is subject to a spectral transformation when biased estimation is used, which has the unfortunate consequence that the computation of the inverse in (25) gives severe numerical difficulties for the specific problem we are dealing with here, i.e., gravity field determination from gradiometry alone, compare Section 4.2 and 4.4.

4 Results

To compute the bias, true coefficients f are needed. For this purpose the spherical harmonic expansion OSU91A, Rapp *et al.* (1991), was taken. All results are relative to GRS80.

Table 2. Regularization with signal constraint

<i>I</i>	obs.	α	$\frac{MSE(i)}{MSE(d)}$	$\frac{\Delta f \Delta f^T}{Q_f}$	$\varepsilon(inv)$
90	<i>yy</i>	1.5	10^5	42	10^{-7}
	<i>zz</i>	3.2	4	10^{-4}	10^{-10}
92.5	<i>yy</i>	0.6	14	1	10^{-8}
95	<i>yy</i>	0.7	16	3	10^{-7}
97.5	<i>yy</i>	1.0	329	91	10^{-7}
	<i>zz</i>	4.9	32	3	10^{-6}
	<i>d</i>	2.5	1	2	10^{-8}
100	<i>yy</i>	520	10^3	10^4	10^{-9}

4.1 Tikhonov regularization

Minimization of the trace of the Mean Square Error Matrix gives the optimal α . For TR this α is found after several iterations. In Table 2 the results for the regularization with the signal constraint are listed. The first two columns define the mission. The third column gives the optimal α for which the trace of the MSEM is minimal. The α 's for the different types of observations cannot be directly compared since the normal matrix N differs from one type of observation to another. Better comparable are the α 's for the same type of observation when only the inclination varies. In general α increases when increasing the size of the polar gap, V_{yy} at 90 degrees is an exception. As already mentioned the zonal harmonics cannot be determined in this case, resulting in a more severe ill-posed problem than e.g. V_{yy} at 92.5 degrees.

In the fourth column the total error of the specific mission is given with respect to the total error of V_d . Only V_{zz} in a polar orbit gives the same level of precision, the other missions are one to five orders worse.

When the trace of the bias part is compared with the trace of the propagated error, column five, one sees that only for V_{zz} , $I = 90$ the bias can be neglected. The bias and the propagation error for V_{yy} at $I = 92.5$ and $I = 95$ and for V_{zz} and V_d at $I = 97.5$ are of the same order of magnitude. For the three remaining missions the bias is two to four orders larger. In fact V_{zz} , $I = 90$ is the only gradient for which the *total* gravity field spectrum can be estimated. The remaining five gradients are not or not sufficiently sensitive to some frequencies of the gravity field, Koop (1993).

The last column lists the maximum error of the inverse, $\varepsilon(inv)$. It was checked by how much

$$(N + \alpha K)^{-1}(N + \alpha K) \quad (26)$$

differed from the identity matrix. The error should ideally be something like 10^{-15} which is the computer round-off error. However, as is evident the maximum error is several orders larger. This means that the stabilized inverse is not as stable as one would like, undermining the value of the computed 'optimal' solution.

4.2 Generalized biased estimation

In Table 3 the results for generalized biased estimation are listed.

Table 3. Generalized biased estimation results

I	obs.	$\frac{MSE(i)}{MSE(d)}$	$\frac{\Delta f \Delta f^T}{Q_f}$	$\frac{MSE(qbe)}{MSE(tr)}$	$\varepsilon(inv)$
90	yy	10^4	24	0.7	10^{-7}
	zz	1	10^{-4}	1	10^{-10}
92.5	yy	273	169	72	10^{-6}
95	yy	10^3	10^3	10^3	10^{-5}
97.5	yy	10^4	10^3	97	10^{-3}
	zz	10^4	15	10^3	10^{-2}
	d	1	2	4	10^{-7}
100	yy	10^4	10^3	34	10^{-3}

Again one sees that only the total error of V_{zz} in a polar orbit is comparable with that of V_d , the other missions are much worse. Apart from the $V_{zz}, I = 90$ case the bias cannot be neglected.

When the total error of the two types of stabilization is compared, we see that the performance of generalized biased estimation is worst in most cases, the missions in a polar orbit give about the same error. These results are in conflict with what one would expect. Since GBE uses several regularization parameters and TR only one, GBE is expected to produce a smaller error, if the parameters are tuned properly. In Section 4.4 this is studied in more detail and an explanation is given.

The maximum error of the inverse is even larger for GBE, the computed inverse is probably unrealistic with the exception of V_{yy} and V_{zz} at $I = 90$ and $V_d, I = 97.5$.

4.3 Error propagation for selected missions to geoid heights

A further understanding of the accuracy of the solutions can be obtained by propagating the errors in \bar{C}_{lm} and \bar{S}_{lm} , as described by the MSEM, to e.g. geoid heights. Because of the mission design (circular orbit, exact repeat, no data gaps) and because the \bar{C}_{lm} and \bar{S}_{lm} errors are assumed to be equal, the Q_f error propagation to geoid heights becomes longitude independent. The normal matrix is block diagonal, orders are independent, even and odd orders are separated, Koop (1993). Furthermore we took the bias part to be of the same block diagonal structure, hence the propagation of the MSEM to geoid heights is longitude independent as well, compare Figure 1. Taking only the block diagonal part of the MSEM as described means neglecting the correlation between orders and between even and odd degrees of the same order. Unfortunately we are forced to do so since the computation of $\Delta f \Delta f^T$ results in a full, $O(L^2 \times L^2)$, matrix, which cannot be handled on nowadays computers.

In Figure 1 the geoid height errors for $V_d, I = 97.5$ are shown. In the polar regions GBE becomes very unrealistic, in the area with observations both methods result in the same geoid error.

Comparing TR at the same inclination $I = 97.5$ for all three types of observables, V_d, V_{yy}, V_{zz} , one sees that in the area with observations V_d and V_{zz} give the same result, whereas V_{yy} is much worse, Figure 2. The GBE errors are not propagated to geoid heights for all these cases because the inverse computation is too unrealistic, compare Table 3.

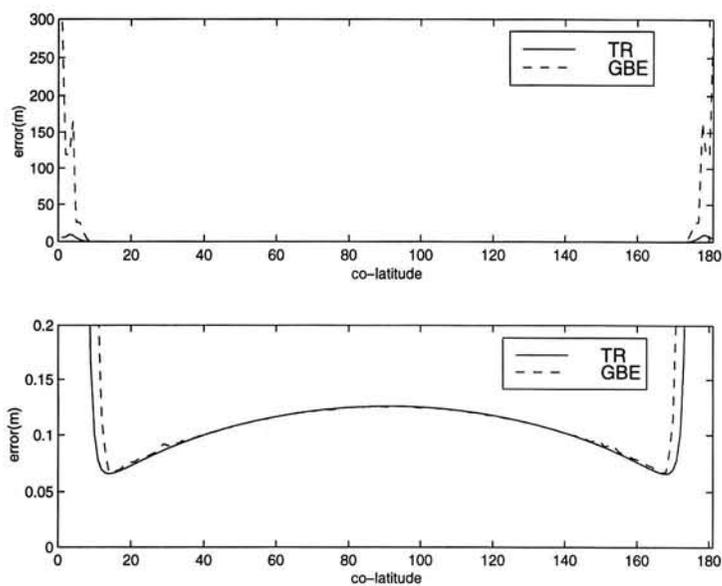


Fig. 1. MSE propagated to geoid heights for V_d , full error (top) and zoom-in (bottom).

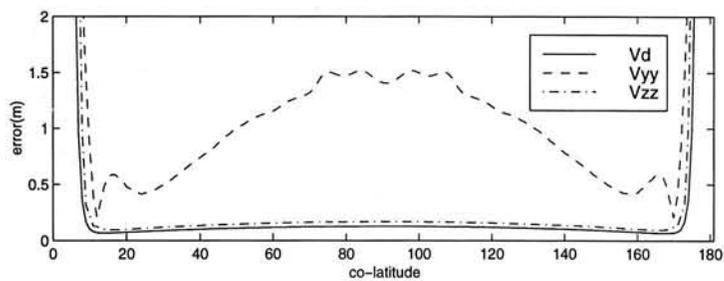


Fig. 2. MSE propagated to geoid heights for V_d, V_{yy}, V_{zz} .

Table 4. Regularization with signal constraint, $m \geq 10$

I	obs.	α	$\frac{MSE(i)}{MSE(d)}$	$\frac{\Delta f \Delta f^T}{Q_f}$	$\varepsilon(inv)$
90	yy	0.7	210	0.4	10^{-11}
	zz	0.8	2	10^{-4}	10^{-12}
92.5	yy	0.6	158	0.3	10^{-11}
95	yy	0.6	98	0.2	10^{-12}
97.5	yy	0.6	54	0.1	10^{-12}
	zz	0.4	3	0.4	10^{-11}
	d	0.5	1	10^{-4}	10^{-13}
100	yy	1.0	84	1.3	10^{-11}

4.4 Exclusion of long wavelengths

The above problems (Section 4.2) are entirely due to the polar gaps. The low orders cannot be determined well because of these gaps. GBE is affected more, which is inherent to the solution. To the small singular values σ_i^2

$$d_i = \langle f, v_i \rangle^{-2} \quad (27)$$

is added. However, the small singular values correspond for a polar gap not only to high frequencies but also to low orders. Consequently, the product $\langle f, v_i \rangle$ becomes large since low frequencies have high energy. Therefore, the squared inverse becomes small which means that there is hardly any stabilization for the small singular values involved with low orders.

This explanation has been tested by excluding the orders $m = 0 - 9$, i.e., we only solve for the spherical harmonic coefficients of degree and order between 10 and 180. The results are given below.

4.4.1 Tikhonov regularization

The results for stabilization by constraining the signal, with $m \geq 10$, are listed in Table 4.

The optimal α 's are in general smaller compared to the full solution and have about the same size. This is no surprise since the major part of the badly solved coefficients has been removed. Therefore, less stabilization is required. Note, however, that a direct comparison of Table 2 and 4 should be done with some care because the normal matrices changed.

When the missions are compared with $V_d, I = 97.5$, one sees that V_{yy} performs better for increasing polar gap. Supposedly this is a consequence of the change of direction of the measurements with respect to an earth fixed reference frame. Both V_{zz} missions give about the same quality as the V_d mission.

The bias is negligible for $V_{zz}, I = 90$ and $V_d, I = 97.5$. It reaches up to 40% for most of the other missions with a minimum of 10%. The bias is 1.3 times larger than the propagated error for $V_{yy}, I = 100$, probably because orders above $m = 9$ are affected here by the polar gap, compare the rule of thumb relating the polar gap and the affected orders in van Gelderen and Koop (1997).

The maximum error in the inverse computation is small, less than 10^{-11} . One can therefore conclude that indeed optimal solutions have been found.

Table 5. Generalized biased estimation results, $m \geq 10$

I	obs.	$\frac{MSE(i)}{MSE(d)}$	$\frac{\Delta f \Delta f^T}{Q_f}$	$\frac{MSE(gbe)}{MSE(tr)}$	$\varepsilon(inv)$
90	yy	169	0.4	0.8	10^{-11}
	zz	2	0.01	1.0	10^{-12}
92.5	yy	124	0.3	0.8	10^{-12}
95	yy	74	0.2	0.8	10^{-12}
97.5	yy	46	0.1	0.8	10^{-12}
	zz	3	0.1	0.9	10^{-11}
	d	1	0.01	1.0	10^{-13}
100	yy	31	1.6	1.0	10^{-10}

4.4.2 Biased estimation

The results for generalized biased estimation, with $m \geq 10$, are listed in Table 5.

Again the V_{zz} missions give about the same MSE as the V_d mission. Here also V_{yy} performs better for increasing polar gap. The V_{yy} missions have a mean square error one to two orders larger than the reference MSE.

The bias for V_{zz} in a polar orbit and V_d at 97.5 degrees is 100 times larger than for TR but still very small: 1% of the propagated error. The relative bias for the remaining missions is about the same for GBE and TR when $m \geq 10$.

The mean square error of GBE now is equal to that of TR or slightly smaller. Again the errors in the inverse show a dramatic improvement with respect to the results in Table 3. The maximum error is less than 10^{-10} .

5 Conclusions

Generalized biased estimation is not suited to stabilize an ill-posed problem when part of the long wavelength signal (high energy) is causing the instability. Tikhonov regularization is better applicable, although the errors of the inverse computation are too large. In general: the method of regularization one has to choose depends on the cause of the instability.

Not estimating the low degree and order harmonic coefficients more or less solves this problem. However, then one has to rely on a priori coefficients which is not desirable for a dedicated gravity field mission like satellite gradiometry. Other constraints, for example by GPS tracking, are therefore needed.

In most cases the bias cannot be neglected. Propagated to geoid errors these biases, and the total error, become largest in the polar areas where no measurements are available. The errors for GBE are not propagated to geoid heights when the error in the inverse is too large. Even when the low degree and orders are excluded the bias can reach up to 40% of the propagated error.

The regularization parameter α should be determined with care since large variations exist in the α giving the minimum Mean Square Error. It is therefore not justified to just pick any number one likes. The α 's shown here are not realistic in the sense that no real or simulated

observations have been used. Future research should use observations and compare the results with those reported here.

Encouraging is the finding that for Tikhonov regularization, excluding the actual polar gaps, and assuming that not only V_{yy} observations are available, the quantity of greatest interest in the space domain, the geoid undulation, can be precisely obtained.

Acknowledgement The computations were partially performed in C++ which was facilitated by the matrix library `newmat08` developed by R. Davies.

References

- Blaser, J., Cornelisse, J., Cruise, A., Damour, T., Hechler, F., Hechler, M., Jafry, Y., Kent, B., Lockerbie, N., Paik, H., Ravex, A., Reinhard, R., Rummel, R., Speake, C., Sumner, T., Touboul, P., and Vitale, S. (1996). STEP Satellite Test of the Equivalence Principle. Report on the Phase A study. ESA SCI(96)5.
- Engl, H., Hanke, M., and Neubauer, A. (1996). *Regularization of Inverse Problems*. Kluwer Academic Publishers.
- ESA (1996). Gravity Field and Steady-State Ocean Circulation Mission. Report for assessment. ESA SP-1196(1).
- Koop, R. (1993). Global gravity field modelling using satellite gravity gradiometry. Publications on geodesy. New series no. 38, Netherlands Geodetic Commission.
- Louis, A. (1989). *Inverse und schlecht gestellte Probleme*. Teubner.
- Rapp, R., Wang, Y., and Pavlis, N. (1991). The Ohio State 1991 geopotential and sea surface topography harmonic coefficient models. Report No. 410, Ohio State University.
- Rummel, R., van Gelderen, M., Koop, R., Schrama, E., Sansò, F., Brovelli, M., Migliaccio, F., and Sacerdote, F. (1993). Spherical harmonic analysis of satellite gradiometry. Publications on geodesy. New series no. 39, Netherlands Geodetic Commission.
- Schrama, E. (1990). Gravity field error analysis: application of GPS receivers and gradiometers on low orbiting platforms. TM 100769, NASA.
- Tikhonov, A. and Arsenin, V. (1977). *Solutions of ill-posed problems*. Winston and Sons.
- van Gelderen, M. and Koop, R. (1997). The use of degree variances in satellite gradiometry. *Journal of Geodesy*, **71**, 337–343.
- Vinod, H. and Ullah, A. (1981). *Recent advances in regression methods*. Marcel Dekker.
- Xu, P. (1992a). Determination of surface gravity anomalies using gradiometric observables. *Geophysical Journal International*, **110**, 321–332.
- Xu, P. (1992b). The value of minimum norm estimation of geopotential fields. *Geophysical Journal International*, **111**, 170–178.
- Xu, P. and Rummel, R. (1994). Generalized ridge regression with applications in determination of potential fields. *Manuscripta Geodetica*, **20**, 8–20.

Quality Assessment of Geopotential Models by Means of Redundancy Decomposition?

Johannes Bouman

Abstract

The determination of a model of the earth's gravitational potential from satellite observations is an ill-posed problem in the sense that a small change in the data may result in a large change in the solution. A stable solution is obtained by adding a priori information to the system of normal equations. One way to describe the quality of the solution is to assess how much the observations contribute to the solution and how much the a priori information. The redundancy number, which is associated with internal reliability, is sometimes used as a measure of this contribution. It is shown here that this is not strictly correct and an alternative method is developed. Moreover, a contribution measure for biased estimators is given.

1 Introduction

Suppose one wants to determine a model of the earth's gravitational field from satellite methods, e.g. satellite tracking, satellite gravity gradiometry, etc. It is then well known that the determination is an inverse or ill-posed problem requiring regularization or stabilization, Rummel *et al.* (1979). In Geodesy such a stabilization is often looked upon as collocation, for example Marsh *et al.* (1988), and the estimated parameters are assumed to be unbiased. Another approach comes from Tikhonov regularization which stabilizes the unstable problem in the same way as collocation, however, inherent to the method is the *regularization error* in the estimated parameters, Louis (1989). It can be shown that this regularization error is equal to the bias in biased estimation, Bouman and Koop (1997).

If one has computed a geopotential model from satellite observations, this model has little value unless one knows its *quality*. But how should 'quality' be described, does there exist a uniform measure? The answer is probably 'no'. First of all a quality assessment can only start when we agree on the (un)biasedness of the solution. Secondly, if we do agree on this issue there are a number of possibilities to gain insight in the 'overall quality'. One could, for example, use the mean square error, or one could propagate the errors in the geopotential model to physical quantities of interest like geoid heights. Other measures are the signal-to-noise ratio of the parameters solved for, the correlation between parameters, reliability, etc. The quality assessment method considered here is *redundancy decomposition* as described in Schwintzer (1990). Basically the redundancy number is a relative comparison of the a priori and

a posteriori error variances of the measurements and can be associated with internal reliability.¹ The smaller the redundancy part of an observation, the larger an error in that observation must be in order to be detectable, Bouman (1993).

The outline of this paper is as follows. Firstly, the linear model relating the measurements to the unknowns is described as well as the method to compute a stable solution. Then Schwintzer's method of quality assessment is recalled and our objections against his interpretation of the redundancy number. We propose a different quality assessment based on the gain matrix from Kalman filtering and generalize the method for biased estimators. Finally the conclusions are listed.

2 Model description and stabilized solution

2.1 Model description

We assume that the relation between the observations, y , and the unknowns to be solved for, x , is linear

$$y = Ax + e \quad (1)$$

where A is the linear or linearized model and e is the vector of misclosures. The latter contains observation errors and model errors. These errors have zero expectation, i.e.

$$E\{y\} = Ax, \quad D\{y\} = P^{-1} \quad (2)$$

with P^{-1} the error covariance matrix of the observations.

Although it is not important in this paper to specify the exact nature of the observations and the unknowns, one could for example think of lumped coefficients as observations and of spherical harmonic coefficients as unknowns. It is important, however, to note that (1) is an abstract notation of the discretized version of an integral equation of the first kind, which is ill-posed. In satellite geodesy this is, for example, due to the downward continuation term $(r/R)^{n+1}$, with $r > R$, which amplifies the measurement noise.

2.2 Solution with collocation

The least-squares collocation solution of (1) is

$$x_c = (A^T P A + C_{xx}^{-1})^{-1} A^T P y \quad (3)$$

or

$$\begin{aligned} x_c &= C_{xx} A^T (A C_{xx} A^T + Q_y)^{-1} y \\ &= C_{xy} (C_{yy} + Q_y)^{-1} y \end{aligned} \quad (4)$$

where C_{ij} are signal covariance matrices between i and j and $Q_y = P^{-1}$ is the measurement error covariance matrix. The equality of (3) and (4) is for example derived in Rummel *et al.* (1979); Bouman (1993). The least-squares collocation solution x_c minimizes

$$J(x) = \|Ax - y\|_P^2 + \|x\|_{C_{xx}}^2.$$

¹The a posteriori errors or residuals are defined as the differences between the real and computed observations after least-squares adjustment, whereas the a priori errors can only be estimated, for example on basis of knowledge of the involved measurement device.

Here, both noise and signal are minimized. C_{xx} could for example be Kaula's rule.

One of the assumptions in collocation is that the signal is on the average equal to zero, Moritz (1980):

$$M\{x\} = 0$$

with M the average over the whole sphere. If the disturbing potential is defined in the proper manner this is true, however, it is not true in general, for example when individual potential coefficients have to be estimated.

Solution (3) also solves (2) extended with zero observations for all unknowns

$$\begin{aligned} E\left\{\begin{pmatrix} y \\ z \end{pmatrix}\right\} &= \begin{pmatrix} A \\ I \end{pmatrix} x, \\ D\left\{\begin{pmatrix} y \\ z \end{pmatrix}\right\} &= \begin{pmatrix} P^{-1} & 0 \\ 0 & C_{xx} \end{pmatrix} \end{aligned} \quad (5)$$

with $z = 0$, and this will be the point of departure for the redundancy decomposition. The error covariance matrix of the estimated parameters is, see (1) and (5)

$$Q_x = (A^T P A + C_{xx}^{-1})^{-1} \quad (6)$$

which can be used in error propagation.

2.3 Solution with Tikhonov regularization

The least-squares solution of (1) can be obtained by minimizing

$$J(x) = \|Ax - y\|_P^2. \quad (7)$$

Dealing with ill-posed problems, however, one does not get a stable solution by minimizing (7). Imposing an additional constraint on x does give a stable solution. Minimizing

$$J_\alpha(x) = \|Ax - y\|_P^2 + \alpha \|x\|_{C_{xx}^{-1}}^2$$

results in

$$x_t = (A^T P A + \alpha C_{xx}^{-1})^{-1} A^T P y = A_\alpha^+ y \quad (8)$$

where $\alpha > 0$ is the regularization parameter.

Of course, there seems to be no essential difference between equations (3) and (8) apart from α . Note, however, that in deriving Tikhonov regularization we made no assumptions about x other than that its total energy is bounded. Thus (8) gives a biased solution unless the average of x is indeed zero. The difference between the exact solution x_e from exact data y_e and the regularized solution x_t is

$$x_t - x_e = A_\alpha^+ (y - y_e) + (A_\alpha^+ - A^+) y_e \quad (9)$$

with $x_e = A^+ y_e$ and A^+ is the generalized inverse of A . The first term on the right-hand side in (9) is called the data error, the second the regularization error Louis (1989). The latter can be shown to be equal to the bias in biased estimation, Bouman and Koop (1997). For biased estimation compare Vinod and Ullah (1981); Xu (1992a, 1992b).

Note that the extended model (5) is not valid here since $E\{0\} = x$ or $E\{x\} = 0$ doesn't necessarily make sense.

3 Contribution of the observations to the collocation solution

Schwintzer (1990) uses the redundancy number as a measure of the contribution of the observations to the solution of the unknowns. Therefore, the redundancy number is discussed first, followed by Schwintzer's interpretation of it. Finally, we present an alternative measure because we do not agree with this interpretation. For further discussion it is referred to Bouman (1993).

3.1 Redundancy number

Consider the linear relationship $E\{y\} = Ax$, where the number of observations is m and the number of unknowns is n , $m \geq n$. Further, let the error covariance matrix $P^{-1} = Q_y$ be diagonal, i.e. the measurements are uncorrelated. The redundancy is $r = m - n$ and it can be shown that (Teunissen, 1994, p. 49)

$$E\{\hat{e}^T Q_y^{-1} \hat{e}\} = \text{trace}(Q_{\hat{e}} Q_y^{-1})$$

and because (Teunissen, 1994, p. 55)

$$E\{\hat{e}^T Q_y^{-1} \hat{e}\} = m - n$$

it holds that

$$\text{trace}(Q_{\hat{e}} Q_y^{-1}) = m - n = r \quad (10)$$

with Q_y the error covariance matrix of y , $Q_{\hat{e}}$ the covariance matrix of \hat{e} and \hat{e} the vector minimizing $e^T Q_y^{-1} e$, $e = y - Ax$. The least-squares solution of x is \hat{x} , and $\hat{y} = A\hat{x}$, $\hat{e} = y - \hat{y}$.

The elements on the diagonal in (10) are denoted as $r_i : [Q_{\hat{e}} Q_y^{-1}]_{ii} = r_i$. The sum of all r_i is

$$\sum_{i=1}^m r_i = r,$$

r_i is the i -th local redundancy number. It is a measure of the extent with which the observation y_i contributes to the total redundancy. Because $Q_{\hat{e}} = Q_y - Q_{\hat{y}}$ (Teunissen, 1994, p. 60) we can write

$$r_i = [(Q_y - Q_{\hat{y}}) Q_y^{-1}]_{ii} = [I - Q_{\hat{y}} Q_y^{-1}]_{ii} = 1 - \frac{\sigma_{\hat{y}_i}^2}{\sigma_{y_i}^2}.$$

Therefore, $0 \leq r_i \leq 1$ since $0 \leq \sigma_{\hat{y}_i}^2 \leq \sigma_{y_i}^2$ (if not the error variances $\sigma_{\hat{e}_i}^2$ could become negative).

Internal reliability. The r_i can be associated with the internal reliability, Förstner (1979a, 1979b); Teunissen (1995), which is a measure of the model error that can be detected with a certain probability γ_0 , for example $\gamma_0 = 80\%$. The minimal detectable bias (mdb) of an observation y_i is, Teunissen (1995):

$$|\nabla_i| = \sigma_{y_i} \sqrt{\frac{\lambda_0}{r_i}}$$

with the noncentrality parameter λ_0 which depends on the choice of γ_0 . The mdb tells us that an error of size $|\nabla_i|$ in observation y_i can be detected with a probability of γ_0 , the power of the test. Therefore, the smaller the redundancy part of an observation, r_i is small, the larger an error in that observation must be in order to be detectable.

3.2 Schwintzer's interpretation

Consider the extended model (5). The redundancy number of the zero observation z_i is

$$r_{z_i} = 1 - \frac{\sigma_{\hat{z}_i}^2}{\sigma_{z_i}^2} \quad (11)$$

where $\sigma_{\hat{z}_i}^2$ and $\sigma_{z_i}^2$ are the i -th diagonal elements of $Q_{\hat{z}}$ and Q_z respectively. Here we have $Q_z = C_{xx}$ and $Q_{\hat{z}} = Q_x$ since $\hat{z} = Ix_c$, see also (6). Schwintzer (1990) uses the local redundancy number (11) as a measure for the contribution to the solution from the observations y . He states: "The partial redundancy r_{z_i} , (...), reflects the contribution of the a priori information to the corresponding results for C_{lm} or S_{lm} in relation to the contribution coming from the real data." (C_{lm} and S_{lm} are coefficients of a spherical harmonic series and the a priori information are the zero observations with weight matrix C_{xx}^{-1} .)

There is some truth in this. Specifically, the zero observations are uncorrelated and $E\{z\} = Ix$. Any redundancy of an observation z_i , that is any verifiability of z_i , is due to the observations y . If r_{z_i} is close to one or $r_{z_i} = 1$, then the corresponding zero observation has excellent internal reliability. Because of the uncorrelated zero observations, the redundant part of z_i has to come from the 'real' observations y and these observations contribute 100% to the verification of z_i . However, we had $E\{z_i\} = x_i$ and therefore one could say that y contributes 100% to the solution of x_i . On the other hand, if $r_{z_i} = 0$ then the corresponding zero observation has poor reliability and y does not contribute to the solution of $z_i = x_i$.

However, our objection to this line of reasoning is that the redundancy has to do with the reliability of observations and it has nothing to do with the unknowns. Moreover, when the a priori information are for example correlated coefficients of an earlier solution, it is not obvious how to explain the redundancy number. Finally note that the expression of the mdb in terms of r_i is derived under the assumption that Q_y is diagonal, Teunissen (1995).

3.3 Alternative contribution measure

Instead of the redundancy number one can also look at the change in the error covariance matrix of the unknowns due to the adding of the zero observations. The contribution of the observations to the solution of the unknowns is defined as follows:

$$contr_y = Q_x Q_{x,y}^{-1} = (A^T P A + C_{xx}^{-1})^{-1} A^T P A \quad (12)$$

where $Q_{x,y}^{-1}$ is the least-squares normal matrix and P not necessarily diagonal. This comparison makes sense. The larger the weight of the observations relatively to the prior information, the larger $contr_y$ gets. Conversely, the smaller the weight matrix P is with respect to C_{xx}^{-1} , the smaller $contr_y$ gets. Of course also the design matrix A is important in the comparison of the observations with the a priori information, since different A 's may give different relative influence of C_{xx}^{-1} . However, we assume A as being a given constant matrix. The contribution of y to an unknown x_i corresponds to the i -th diagonal element of the matrix $contr_y$.

The contribution of the zero observations to the solution of the unknowns is analogously

$$contr_z = Q_x Q_{x,z}^{-1} = (A^T P A + C_{xx}^{-1})^{-1} C_{xx}^{-1} \quad (13)$$

Note that the sum of the two contributions equals I_n , which indicates that the total information for one estimator comes from the observations and the a priori information together. Further

note that the i -th diagonal element of $contr_z$ is

$$contr_{z_i} = \frac{\sigma_{x_i}^2}{\sigma_{x,z_i}^2} = \frac{\sigma_{z_i}^2}{\sigma_{z_i}^2}$$

since $C_{xx} = Q_{x,z}$ is a diagonal matrix. Therefore,

$$contr_{y_i} = 1 - contr_{z_i} = 1 - \frac{\sigma_{z_i}^2}{\sigma_{z_i}^2}$$

which happens to be equal to (11). On purpose we use 'happens' since in case of correlated a priori information these two numbers will not be equal. Still (12) can easily be used and makes sense.

As is shown in Bouman (1993), equation (13) equals the gain matrix in Kalman filtering which describes the precision improvement of the estimated unknowns when additional measurements become available, Salzmann (1993). It is, however, not allowed to speak about precision improvement. The least-squares solution from y alone is unstable and has to be avoided. Equation (12) is just a measure for the contribution of y to x_i .

4 Contribution of the observations to the biased solution

The above derivations are all based on the assumption of unbiasedness of the estimator. However, the solution might be biased and the precision of the solution can no longer be described with the propagated observation error alone: the bias has to be included as well.

4.1 The mean square error matrix

Remark: The equations involved with the mean square error matrix (and bias) become more complex compared to the unbiased case. Therefore, we will abandon the weight matrices P and C_{xx} for the moment, i.e. $P = C_{xx} = I$.

The propagated data error is from (8)

$$Q_x = (A^T A + \alpha I)^{-1} A^T A (A^T A + \alpha I)^{-1}$$

where explicitly $E\{z\} \neq x$ has been used. The bias in the solution is

$$\Delta x = -(A^T A + \alpha I)^{-1} \alpha x$$

which can be derived from (9), see for example Bouman and Koop (1997). The total error or Mean Square Error Matrix (MSEM) consists of the propagated error and the bias

$$MSEM = Q_x + \Delta x \Delta x^T = (A^T A + \alpha I)^{-1} (A^T A + \alpha^2 x x^T) (A^T A + \alpha I)^{-1}.$$

The contribution of y now is, compare with (12)

$$contr_y = MSEM \times MSEM^{-1} \Big|_{\alpha=0} = (A^T A + \alpha I)^{-1} (A^T A + \alpha^2 x x^T) (A^T A + \alpha I)^{-1} A^T A. \quad (14)$$

Remark: Including weight matrices (14) becomes

$$\text{contr}_y = (A^T P A + \alpha C_{xx}^{-1})^{-1} (A^T P A + \alpha^2 C_{xx}^{-1} x x^T C_{xx}^{-1}) \times (A^T P A + \alpha C_{xx}^{-1})^{-1} A^T P A. \quad (15)$$

Now suppose that C_{xx} can be represented by Kaula's rule and that x are potential coefficients. Then $C_{xx} \approx \text{diag}(x x^T)$ since Kaula's rule approximates the true values, or

$$\text{contr}_y \approx (A^T P A + \alpha C_{xx}^{-1})^{-1} (A^T P A + \alpha^2 C_{xx}^{-1}) (A^T P A + \alpha C_{xx}^{-1})^{-1} A^T P A \stackrel{\alpha \equiv 1}{=} (A^T P A + C_{xx}^{-1})^{-1} A^T P A,$$

neglecting the off-diagonal terms of $x x^T$ and for $\alpha = 1$. Hence, for this special case the MSEM is equal to (6) and (15) equals (12).

4.2 The mean square error

An alternative to the use of the MSEM is the spectral decomposition of the Mean Square Error (MSE). The trace of the MSEM is defined as the MSE, which can be shown to be, e.g. Bouman (1993):

$$\text{MSE} = \sum_{i=1}^n \frac{\sigma_i^2 + \alpha^2 \langle x, v_i \rangle^2}{(\sigma_i^2 + \alpha)^2}$$

with the singular value decomposition of $A = U \Sigma V^T$ and v_i is a column vector of V and σ_i is the i -th diagonal element of Σ . Comparing for a single i the MSE with and without α gives

$$\text{contr}_{y_i} = \frac{\text{MSE}_i}{\text{MSE}_i|_{\alpha=0}} = \frac{\sigma_i^2 + \alpha^2 \langle x, v_i \rangle^2}{(\sigma_i^2 + \alpha)^2} \sigma_i^2. \quad (16)$$

When the eigenvalue σ_i^2 is large (compared to α) it means that the unknown x_i is represented well by the measurements and $\text{contr}_{y_i} \approx 1$. On the other hand, when σ_i^2 is small compared to α , $\text{contr}_{y_i} \approx 0$, as required. Of course also $\langle x, v_i \rangle$ has a certain size and cannot be neglected. However, in the time domain $\langle x, v_i \rangle$ has comparable size for various i since v_i is an orthonormal vector and the energy of the signal is expected to be about equal at different locations. In the frequency domain $\langle x, v_i \rangle$ is large when the long wavelengths are poorly determined and will be smaller for shorter wavelengths (less energy). Thus, only poorly determined long wavelengths could cause problems because the numerator becomes larger. Then the above contribution measure probably fails.

4.3 Discussion on the mean square errors

Comparison of MSEM and MSE. Let the singular value decomposition of A be $A = U \Sigma V^T$, with U and V orthonormal matrices, i.e., $U^T = U^{-1}$, $V^T = V^{-1}$ and Σ is the matrix with the singular values on the main diagonal. Some elementary calculation then gives

$$\text{contr}_y = V(\Sigma^2 + \alpha I)^{-1} (\Sigma^2 + V^T \alpha^2 x x^T V) \times (\Sigma^2 + \alpha I)^{-1} \Sigma^2 V^T.$$

The trace of contr_y now yields

$$\sum_{i=1}^n \text{contr}_{y_i}$$

where contr_{y_i} is defined in (16), see also Xu and Rummel (1994). However, the diagonal elements of contr_y are not the same as contr_{y_i} , only their sum consists of these elements. The elements $[\text{contr}_y]_{ii}$ are build of contr_{y_i} subject to a 'spectral transformation' by V and V^T . At first sight we prefer the measure (15) since it is directly comparable with (12). The measure based on MSE maybe usefull as well, numerical examples could further clarify this issue.

Computational aspects. Assuming that x is known (which it is not since we are trying to estimate x), a severe problem in the computation of MSEM and $contr_y$ is that xx^T gives a full matrix. When for example potential coefficients are the unknowns up to degree and order 180, this matrix is too large to be handled on nowadays computers. One could take the block diagonal part of the full matrix but the neglected parts may be of the same order as the block diagonal elements. Unfortunately we do not know of any better alternative yet.

The realization of the MSE requires a singular value decomposition or, better, an eigenvalue decomposition of the positive definite normal matrix $A^T A$. The eigenvalue decomposition is better since it is a less expensive decomposition in terms of computation time. However, if the normal matrix does not have a special structure, also the eigenvalue decomposition may give trouble. Again, one could take only the block diagonal part for example.

Weight matrices. The weight matrices P and C_{xx} are in general not equal to the identity matrix of course. Therefore, it may seem necessary to redo the above derivations for the more general case. Fortunately this is not necessary, since a transformation to the standard form: minimize

$$J_\alpha(x) = \|Ax - y\|^2 + \alpha\|x\|^2$$

is always possible, see e.g. Hansen (1997); Bouman (1993).

How to compute the bias. It was noted above that we do not know x , hence, we cannot compute the bias or the mean square error. Instead, we could use the biased estimator x_l , but this might underestimate the bias, Xu (1992b).

Specifically for satellite geodesy one might choose some existing gravity model like OSU91A, Rapp *et al.* (1991), to compute the bias or one can take the size of the coefficients according to Kaula's rule and the sign following from the satellite solution, Koop (1993). Then the bias will not be underestimated and the error assessment is probably more reliable.

5 Conclusions

The answer to the question: 'Quality assessment of geopotential models by means of redundancy decomposition?', must be 'No, unless the a priori information has the simple form (5) and the estimator is unbiased.' Otherwise it is better to use (12) which is related to the gain matrix in Kalman filtering, or (15) when the estimator is biased.

A further understanding of the respective measures and a comparison of them, would be by numerical examples. The major problem for the 'biased' measures is the estimation of the bias and the handling of the large matrix xx^T .

References

- Bouman, J. (1993). *The normal matrix in gravity field determination with satellite methods; its stabilization, its information content and its use in error propagation.* Master's thesis, Delft University of Technology.
- Bouman, J. and Koop, R. (1997). Quality differences between Tikhonov regularization and generalized biased estimation in gradiometric analysis. *DEOS Progress Letters*, **97.1**, 42-48.
- Förstner, W. (1979a). Das Programm TRINA zur Ausgleichung und Gütebeurteilung geodätischer Lagenetze. *Zeitschrift für Vermessungswesen*, **104**, 61-72.
- Förstner, W. (1979b). Das Rechenprogramm TRINA für geodätische Lagenetze in der Landesvermessung. *Nachr. a. d. öff. Vermessungsw. NW*, **12**, 125-166.

- Hansen, P. (1997). *Regularization Tools, A Matlab package for analysis and solution of discrete ill-posed problems, Version 2.1 for Matlab 5.0*. Department of Mathematical Modelling, Technical University of Denmark. <http://www.imm.dtu.dk/~pch>.
- Koop, R. (1993). Global gravity field modelling using satellite gravity gradiometry. Publications on geodesy. New series no. 38, Netherlands Geodetic Commission.
- Louis, A. (1989). *Inverse und schlecht gestellte Probleme*. Teubner.
- Marsh, J., Lerch, F., Putney, B., Christodoulidis, D., Smith, D., Felsenreger, T., Sanchez, B., Klosko, S., Pavlis, E., Martin, T., Williamson, J. R. R., Colombo, O., Rowlands, D., Eddy, W., Chandler, N., Rachlin, K., Patel, G., Bhati, S., and Chinn, D. (1988). A new gravitational model for the earth from satellite tracking data: GEM-T1. *Journal of Geophysical Research*, **93**(B6), 6169–6215.
- Moritz, H. (1980). *Advanced physical geodesy*. Wichmann.
- Rapp, R., Wang, Y., and Pavlis, N. (1991). The Ohio State 1991 geopotential and sea surface topography harmonic coefficient models. Report No. 410, Ohio State University.
- Rummel, R., Schwarz, K., and Gerstl, M. (1979). Least squares collocation and regularization. *Bulletin Géodésique*, **53**, 343–361.
- Salzmann, M. (1993). Least squares filtering and testing for geodetic navigation applications. Publications on geodesy. New series no. 37, Netherlands Geodetic Commission.
- Schwintzer, P. (1990). Sensitivity analysis in least squares gravity field modelling by means of redundancy decomposition of stochastic a priori information. Deutsches Geodätisches Forschungs-Institut, internal report.
- Teunissen, P. (1994). *Mathematische geodesie I, inleiding vereffeningstheorie*. Lecture notes, Delft University of Technology, Faculty of Geodetic Engineering, Delft. (Introduction to adjustment, in english).
- Teunissen, P. (1995). *Mathematische geodesie II, inleiding toetsingstheorie*. Lecture notes, Delft University of Technology, Faculty of Geodetic Engineering, Delft. (Introduction to testing, in english).
- Vinod, H. and Ullah, A. (1981). *Recent advances in regression methods*. Marcel Dekker.
- Xu, P. (1992a). Determination of surface gravity anomalies using gradiometric observables. *Geophysical Journal International*, **110**, 321–332.
- Xu, P. (1992b). The value of minimum norm estimation of geopotential fields. *Geophysical Journal International*, **111**, 170–178.
- Xu, P. and Rummel, R. (1994). Generalized ridge regression with applications in determination of potential fields. *Manuscripta Geodetica*, **20**, 8–20.

Overview of tide gauge systems and averaging techniques

Kyra van Onselen

Abstract

Variations in sea level can be determined using tide gauges, relating the variation in sea level height to a local bench mark. The quality of the resulting time series depends on the measuring accuracy of the tide gauge system used and on the frequency of the measurements. Since sea level variation curves are usually based on monthly or yearly mean sea level heights, the method used to form these mean sea level values also influences the accuracy of the determined curves.

In this article an overview of error characteristics (precision, inherent systematic errors, and limitations) is given for six major tide gauge systems: tide poles, tide poles with float, stilling well tide gauges, reflection tide gauges, subsurface pressure tide gauges, and open sea pressure gauges. In addition a number of methods are discussed which have been used in the past, or are still being used nowadays, to form monthly and annual mean sea level heights.

1 Introduction

Sea level is conventionally monitored using tide gauges, which relate variations in sea level height to a local tide gauge bench mark. Over the years a number of different tide gauge systems have been developed, with varying measuring precision and recording methods. In Section 2 error characteristics will be discussed for the six major tide gauge systems that have been used in the past or are still in use today, i.e., tide poles, tide poles with float, stilling well tide gauges, reflection tide gauges, pressure tide gauges and open sea pressure tide gauges.

The quality of sea level height time series and, consequently, the quality of estimated patterns in sea level variation, not only depends on the measuring accuracy of the tide gauge systems which have been used, but on the measuring frequency as well. In addition, since usually only some kind of mean values (hourly, daily, monthly, or even yearly values) are available, the method used to form these mean values also influences the quality of sea level height time series. These effects of sampling rate and averaging method will be discussed in Section 3.

2 Error characteristics of tide gauge instruments

Local sea level variations have been measured for hundreds of years using a variety of measuring techniques. At first, sea level measuring systems consisted of a vertically mounted pole,

relative to which instantaneous sea level height could be determined by an observer. Since reading sea level height relative to a pole is difficult in the presence of waves, the introduction of a float which moves vertically in a well significantly improved the accuracy of the measurements. This system was further improved by automating the recording of the sea level data, e.g., by means of a pen driven by the float across a chart, which in turn is mounted on a circular drum. According to Pugh (1987), self-recording gauges began operating in the beginning of the nineteenth century. Mechanical recording has a number of disadvantages as well, e.g., distortion of the paper chart due to humidity, friction of the pen, etc. Therefore, digital recording leads to a further improvement in the accuracy of the measurements.

Besides improvements on the "traditional" system of the stilling well with float system, a number of other sea level measuring devices, based on other measuring principles, have been developed during this century. Examples are tide gauges that measure pressure at a fixed point below the sea surface and tide gauges that measure the reflection time between a fixed point and the instantaneous sea surface.

For planning required measurement durations for sites where new tide gauges have to be installed, only the accuracy of state-of-the-art tide gauge systems is important. However, error characteristics of older sea level measuring systems are relevant as well since often older sea level records have to be included in order to estimate specific phenomenon in sea level height variations. For some sites, tide gauge data have been recorded for more than a century, and, in order to determine meaningful results from these data, variations in data quality over the length of the time series (e.g., due to changes in applied measuring techniques) have to be taken into account. Error characteristics and limitations for the six major techniques, i.e., tide poles, tide poles with float, stilling well tide gauges, pressure tide gauges, reflection tide gauges, and open sea pressure tide gauges, will be described in the following sections.

2.1 Tide poles

One of the first systems for monitoring sea level variations consisted of a so-called tide pole, vertically mounted on a site assumed representative for the area of interest. At regular intervals the height of the instantaneous sea level relative to this pole can be read by an observer. The gauge zero should be connected to a permanent bench mark on shore, the so-called tide gauge bench mark.

The precision of sea level measurements from tide pole readings is determined by random reading errors which, according to Montag (1970), depend on the state of the sea surface, the quality of the markers, the distance between the observer and the tide pole, light conditions, and the constitution of the observer. From experiments carried out under different weather conditions, Montag (1970) concluded that the average error due to these reading errors amounts to less than 2 cm. Pugh (1987) is slightly less optimistic, he states a reading precision of 2 cm but only for calm weather conditions; in the presence of waves this precision deteriorates. According to Pugh (1987), experiments show that in the presence of waves with a height of 1.5 meters, experienced observers are able to read the tide pole with a precision of 5 cm.

Apart from these random reading errors, the accuracy of tide pole measurements is also determined by the presence of systematic errors and blunders. Blunders can originate from gross errors either in the readings of the sea level heights or in the registration of these heights. Sys-

tematic errors can result

1. from the construction of the pole itself
e.g., deviations in the scale of the pole, or problems with the connection between the individual 1 meter sections of the pole,
2. from environmental conditions
e.g., due to different illumination during day and night, or
3. they are caused by the operator
e.g., operator has a tendency to read slightly to high values.

Furthermore, sea level in the presence of waves is often determined as the average between the crest and trough of the waves, as a result, due to two different mechanisms waves give rise to additional systematic errors. First of all, as explained by Pugh (1987), for an observer viewing the pole at an angle the apparent trough may actually be the crest of an intermediate wave which is obscuring the true trough at the pole, resulting in an averaged level which is too high. In addition, since water waves differ from a true sine oscillation, averaging between crest and trough values will give a systematic error, which value depends on the state of the sea surface, but on average will be in the order of 1 mm; see Montag (1970) for more details.

Although tide poles have a large number of limitations they have some advantages as well. For example, due to the small amount of technology involved they are relative cheap and easy to install and operate. Consequently, they provide an easy method to check readings of other tide gauge systems, e.g, to check the datum imposed for pressure tide gauge measurements (see Section 2.5).

2.2 Tide pole with float

As indicated by, e.g., Pugh (1987), the problem of reading the tide pole (especially in the presence of waves) can be minimized as follows. A transparent tube is fitted alongside the tide pole, this tube is connected to the sea by means of a narrow inlet tube, which prevents waves from entering the tube connected to the tide pole, see Figure 1. As a result, the height of the sea level inside the tube can easily be read relative to the tide pole.

How well waves are attenuated depends on the so-called time constant of the specific tide pole system. According to Smithson (1997), this time constant is determined by the characteristics of the inlet tube and stilling tube and can be estimated as

$$\tau = \frac{32\nu L_p D_w^2}{g D_p^4} \quad (1)$$

in which g is the gravitational acceleration, ν is the kinematic viscosity of sea water, D_p and L_p are resp. the diameter and length of the inlet tube, and D_w is the internal diameter of the stilling tube.

The amplitude of a wave with frequency ω (in radians) is attenuated by a factor α which can be estimated by

$$\alpha = \frac{1}{\sqrt{1 + \omega^2 \tau^2}} \quad (2)$$

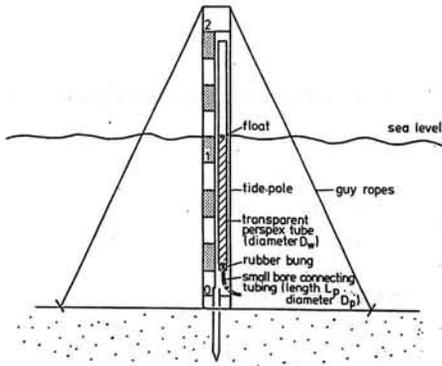


Fig. 1. Tide pole with float; reproduced from Smithson (1997)

In addition, the inlet tube system causes a phase lag (θ) of the wave given by

$$\tan \theta = \omega \tau \quad (3)$$

Both equations are given by Smithson (1997).

A tide pole with a float in a stilling tube has a lot of error characteristics and limitations in common with stilling well tide gauges. Since stilling well tide gauges are one of the more popular tide gauges in use nowadays, these characteristics will be described in detail in the next section. In this section, problems related to using a float and a stilling tube will only be discussed briefly.

The accuracy of tide pole measurements fitted with a float in a stilling tube is determined by mechanical problems and the deviation between the water level in the tube and the open sea level outside. Mechanical problems are, e.g., friction of the float within the stilling tube and deviations in the construction of the pole itself as discussed in Section 2.1.

Deviations between the water level inside and outside the stilling tube depend, e.g., on the state of the sea surface in connection with the size of the inlet and stilling tube. Waves, which cause reading inaccuracies if "normal" tide poles are used, are damped out by the filtering properties of the inlet and stilling tube construction. How well these waves are attenuated depends on the size of these tubes. However, this mechanical filtering of high frequency signals also results in a small time delay of the sea level in the tube as compared to the sea level height variations outside the tube. Since this delay is frequency dependent (see equation 3), this causes an error in measured sea level. As will be discussed in more detail in the following section, an additional problem is that due to environmental conditions (like silting up by sediments and marine growth) the filtering characteristics of the inlet and stilling tube construction might change over time.

Other factors contributing to a deviation between levels inside and outside the stilling tube are differences in salinity and temperature between the "open sea" water and the water inside the tube and pressure deviations due to the inlet tube being placed in a tidal stream. Since these error sources also impede stilling well tide gauges, they will be further discussed in Section 2.3. An additional problem with salinity and temperature differences in the stilling tube is that their effects are extremely difficult to estimate if dye is mixed into the water in order to increase

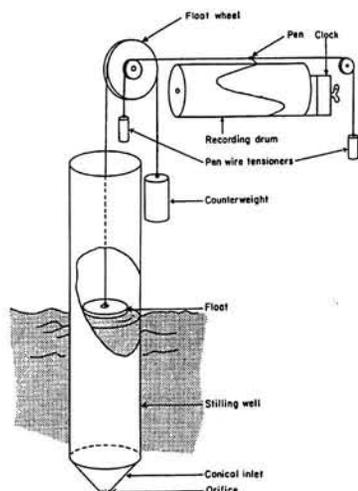


Fig. 2. Stilling well tide gauge with mechanical recording on a drum; reproduced from IOC (1985)

the visibility of the level of water in the tube and/or some substances are added to the water to prevent freezing.

2.3 Stilling well tide gauges

The two previous tide gauge systems have the major disadvantage that tide gauge readings have to be performed by a human observer. Apart from the fact that this facilitates the occurrence of blunders (gross errors in reading or recording of sea level heights) this makes them unsuitable for long term measurement campaigns with a high measuring frequency. This high measuring frequency allows elimination of a.o., high frequency disturbances, and results in more precise (e.g., hourly) mean values.

Since the mid-nineteenth century self-recording gauges have been in operation. They often consist of a so-called stilling well (a tube or other enclosed area, connected to the open sea by an orifice or inlet pipe) with a float. This float is connected to the recording system, which in the past was a mechanical device, whereas newer tide gauges use digital recording. Figure 2 shows an example of a stilling well tide gauge with orifice and mechanical recording on a drum.

In order to prevent aliasing and to allow for the use of a float, high frequency movements of the water surface have to be filtered out mechanically. The stilling well provides a protection of the float system against environmental conditions (e.g., from the effect of wind), whereas the relatively small diameter of the inlet provides a mechanical filter of high frequency movements like waves. Damping characteristics of the stilling well are, mainly, determined by the ratio between the diameter of the inlet and the diameter of the stilling well itself.

Damping effect of the stilling well system is determined by the flow coefficient through the inlet. As explained in detail by, e.g., Sager and Matthäus (1970), this flow coefficient includes the hydraulic losses inherent to the water having to move through the inlet, and has a value between 0 (inlet is completely closed) and 1 (ideal flow conditions). Detailed relations between

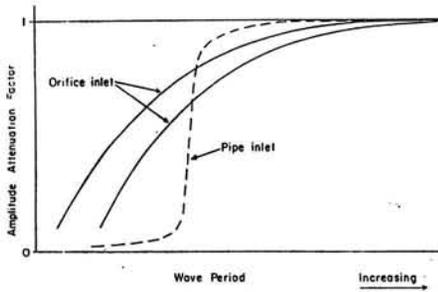


Fig. 3. Relation between attenuation factor and wave period for stilling wells with resp. an orifice and a pipe inlet; reproduced from Smithson (1997)

water level response in the stilling well to resp. constant, linear and periodic variations of the “open” sea water level have been described by Sager and Matthäus (1970), both this is beyond the scope of this article.

The term “inlet” has been used to indicate one of the two types of inlets commonly used, i.e., orifices and inlet pipes. These two types of inlets have different attenuation and (maybe even more important) phase lag characteristics. As explained in detail by Noye (1974), a stilling well using an inlet pipe has a number of advantages compared to a stilling well with an orifice, i.e.,

- using an inlet pipe gives full system response (no attenuation) for a relatively wide range of tidal and long-period oscillations,
- high-frequency phenomena like wind waves do not cause a systematic set-down of the mean water level in the stilling well,
- inlet pipes yield linear systems which allow tidal constituents to be directly corrected for attenuation and phase lag, and
- instead of gradually increasing the attenuation factor with increasing frequency, stilling wells with an inlet pipe have a relatively sharp cut-off.

These effects are illustrated by Figure 3 (reproduced from Smithson (1997)) which shows for both types of inlet the relation between wave period and attenuation factor.

The accuracy of sea level variations based on stilling well tide gauges is determined by the stilling well itself (deviations between the water level inside the stilling well and the “open” sea), by mechanical problems of the float system, and by problems with the recording of data. If digital recording is used to record samples of sea level height at discrete periods in time, based on the mechanical filtering characteristic of the stilling well system, the sampling period has to be carefully selected in order to prevent aliasing of remaining high frequency fluctuations in the stilling well.

Mechanical recording has the advantage that it allows for a continuous monitoring of sea level variations, and thus the problem of aliasing is averted. However, mechanical recording of sea level heights, e.g., on a chart mounted on a circular drum (see Figure 2), yields inaccuracies in recorded sea level heights as well. Main problems, as indicated by e.g., Lennon (1970),

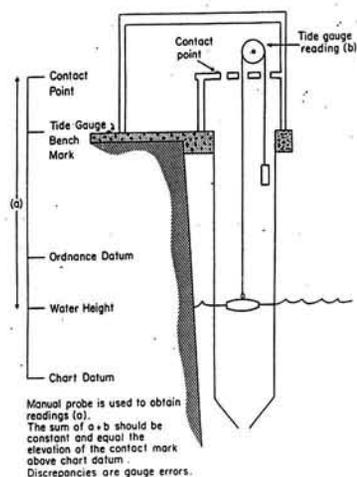


Fig. 4. Datums involved with stilling well tide gauge; reproduced from Smithson (1997)

Van der Made (1987), and Xu (1990), are problems with the time scale (e.g., due to a deviation in the rotational speed of the drum, or even worse, due to variations in this rotation), deformation of the drum, and deformation of the paper (e.g., due to variations in humidity). In addition, as indicated by Pugh (1987), reading of the produced charts is a tedious procedure which is prone to errors. Pugh (1987) estimates that due to the width of the chart trace, precision of sea level height charts is in the order of 2 cm for levels and 2 minutes in time

Mechanical problems with the float system include movements of the float and wire, friction and backlash in float and counterpoise suspension, variable tension in the suspension wire of the float, and variations in the length of the float suspension (again based on Lennon (1970), Van der Made (1987), and Xu (1990)). A number of these mechanical problems can be checked against by the so-called "van de Castelee test". In this test, the distance between the contact point and water level is measured manually (e.g., with a steel tape) and compared to the recorded sea level height, see Figure 4. Plots of the difference between these two distances through time (the so-called "van de Castelee diagrams") give an indication of probable mechanical errors in the tide gauge system. For an "ideally" operating tide gauge the diagram shows a straight line, different types of deviations from these straight line correspond to different types of mechanical problems. For more details the interested reader is referred to, e.g., Smithson (1997).

The last category of limitations on the accuracy of stilling well sea level heights are introduced by the stilling well structure itself. As indicated in the preceding section, mechanical filtering by means of a relative small inlet not only attenuates high frequency signals it also causes a time delay between sea level variations in the "open" sea and resulting water level variations in the stilling well. Furthermore, the water level in the stilling well deviates from that in the "open" sea if differences in density (differences in temperature and salinity) and pressure occur between the water inside and outside the stilling well.

An additional problem with mechanical filtering by means of a relative narrow inlet is that its filtering characteristics change if the diameter of the inlet changes. Environmental causes which can narrow the inlet are siltation, marine growth and accumulation of weed or trash.

Tide gauges situated in some sites will be more prone to clogging up than tide gauges situated in other sites. Besides, as investigated by Cross (1968), the flow through an orifice inlet is asymmetrical, resulting in a pumping-down effect of the water level in the well if significant surging outside the well occurs. Cross (1968) estimates that this effect can give an error of 14 cm in the presence of waves with heights between 1.5 and 3 meters.

Deviations in density (in temperature and salinity) between the sea water and the water inside the stilling well causes different water levels inside and outside the well. As described by Van der Made (1987) if salinity of the "open" sea water changes (e.g., due to variations in runoff of a nearby river), the effect on the density of the water outside the stilling well will be much larger than the effect on the water inside the well. As an example of the impact of this phenomenon on determined sea level heights, Van der Made (1987) shows that for a stilling well with an inlet tube 4 meters below the water level and a density of 1010 kg/m^3 inside the well, an increase in density outside the well from 1000 to 1020 kg/m^3 would yield a deviation between the level inside and outside the well of 4 cm.

Especially stilling well tide gauges at estuary sites can show large deviations between the water level inside and outside the stilling well due to variations in salinity and temperature through the tidal cycle. As explained by, e.g., Pugh (1987), with rising tide the density of estuary water increases. Consequently, at high tide the density of the "open" sea water is higher than inside the stilling well, where the density is an average of the "open" sea density during the filling up of the stilling well. As a result, the water level in the stilling well can be significantly higher than outside the stilling well. In extreme cases, like tide gauges in the river Mersey where the tidal range is in the order of 10 meters, Lennon (1970) estimates that at spring high water the level inside the well can be 6 cm higher than outside the well. Pugh (1987) even gives a difference in water level of 12 cm for a tidal range of 10 meters.

Deviations between the water level inside the stilling well and the "open" sea level can also be caused by differences in pressure. If the stilling well is situated in a tidal stream, the stilling well structure itself causes pressure disturbances by (partly) blocking the flow. In addition, other obstacles in the vicinity of the stilling well can cause pressure differences. The result is a draw-down of the water level, which gives (analogous to differences in density) systematic errors in recorded sea level heights.

Estimates of the precision of stilling well tide gauge measurements and resulting mean hourly values vary widely. For example, Christensen *et al.* (1994) state that tide gauges at Harvest platform estimate sea level heights with a sample standard deviation of ± 1.5 cm around the mean. On the other hand, according to Diamante *et al.* (1987), the standard deviation of a single measurement of sea level height made by a stilling well tide gauge will not be better than about 5 cm. The precision of sea level heights can be improved by forming hourly mean values, since this reduces the influence of high frequency (nearly) random errors as introduced by waves. Hamon and Godfrey (1980) estimate that an adequately maintained tide gauge situated in a suitable location can yield hourly mean values with a precision of 1 cm, whereas daily mean values will have an even higher precision, i.e., a standard deviation of only 1 mm. A suitable location is (probably) a site which is the least possible influenced by local processes, so the tide gauge is e.g., not installed near a large river mouth.

In Section 3 different methods will be discussed which can be used to form longer-term average values like monthly and yearly means. The precision of these mean values will be significantly better than that of the single samples of sea level heights because the effect of most random influences (like the effect of waves) will be strongly reduced. However, a number of

systematic errors (like draw-down by pressure differences and relative high water levels in the stilling well during high tides due to differences in density) also occur in sea level measurements. These errors cannot be reduced by forming monthly or yearly averages of sea level. As explained by Diamante *et al.* (1987) special attention has to be paid to these systematic effects since, if not taken into account by other methods, they will introduce long-period effects which are almost identical to the actual trend in sea level height.

One possible way of dealing with systematic effects is to combine sea level readings for a number of tide gauges in a region. However, it should be noted that some of these systematic effects (like hydrodynamic draw-down) will influence all tide gauges in the same manner and, consequently, it is not possible to completely eliminate these effects by simple averaging tide gauge data over a large region. Combining time series for a number of tide gauges in a region does not necessarily yield a mean time series with an improved accuracy. Although taking an area average will reduce some of the systematic errors present in the individual time series, some new errors will be introduced as well. For example, the accuracy of the combined time series is limited by inaccuracies in the determined height differences between the tide gauge benchmarks of the tide gauges under consideration.

2.4 Reflection tide gauges

As explained in the preceding section, the accuracy of stilling well tide gauge measurements is largely influenced by the filtering characteristics of the (narrow) orifice or inlet pipe (and especially by alterations of the filtering characteristics by, e.g., siltation of the inlet) and by mechanical problems with the float system. According to, e.g., Diamante *et al.* (1987), these systematic errors can be prevented (or at least largely reduced) if this float is replaced by a remote sensor which does not require physical contact with the actual water surface. For this type of sensor it is no longer necessary to attenuate the high-frequency fluctuations which makes mechanical filtering by means of a narrow orifice or inlet pipe superfluous.

In practice, as indicated by Diamante *et al.* (1987), present day acoustic measuring equipment still requires some kind of protective well in order to limit power consumption and procure accurate enough results. However, an open protective well, in which sea level conditions will be more resemblant to the open sea (less sensitive to e.g., density build-ups, etc.), which is less prone to siltation and does not have a strong filtering effect, might be sufficient. On the other hand, for a state-of-the-art station configuration Martin *et al.* (1996) still propose the use of a protective well with an orifice opening. However, they use a 15 cm diameter protective well with a 5-cm diameter orifice, whereas for stilling well gauges Pugh (1987) recommends an orifice to well diameter ratio of 0.1.

Figure 5, which was reproduced from Martin *et al.* (1996), gives a schematic drawing of such a state-of-the-art reflection tide gauge. Sea level height is measured by means of the time taken by an acoustic pulse to travel between the acoustic sensor, the instantaneous sea surface and back. From this two-way travel time and the height of the acoustic sensor relative to the tide gauge bench mark the instantaneous sea level height can be determined.

The two-way travel time (t_r) not only depends on the distance between the acoustic sensor and the sea surface (l_i) but also on the velocity of the acoustic pulse in the sounding tube (C_a), i.e.,

$$t_r = \frac{2l_i}{C_a} \quad (4)$$

For dry air at a temperature of 10° and one atmosphere pressure, sound has a velocity of 337.5

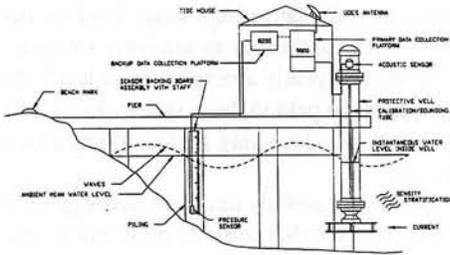


Fig. 5. Reflection tide gauge system; reproduced from Martin *et al.* (1996)

m/s. This implies that, under these conditions, timing must be accurate within $5.9 \cdot 10^{-5}$ seconds, in order to detect a 1 cm change in sea level height.

Velocity of sound in air depends on temperature, pressure and humidity, therefore, in order to accurately determine variations in sea level height, variations in these parameters have to be taken into account. The new reflection tide gauge system as described by Martin *et al.* (1996) consists of a self-calibrating acoustic sensor which compensates for variations in C_a (due to temperature changes) in the sounding tube.

According to Martin *et al.* (1996), sea level heights are determined as average values over 6-minute intervals based on instantaneous measurements with a sampling rate of 1/s. Averages are determined over 3-minute periods, while sample outliers exceeding three times the standard deviation corresponding to this 3 minute interval are removed from the data. The resulting new averages provide one sea level measurement for every 6 minutes with a resolution of ± 1 mm. Smithson (1997) is slightly less optimistic; he estimates that a resolution of only about 3 mm can be achieved by averaging measurements with a sampling rate of 1/s over periods of 3 minutes.

2.5 Subsurface pressure tide gauges

As an other alternative for the stilling well tide gauge, nowadays often subsurface pressure tide gauges are used. These tide gauges measure pressure at a fixed point somewhere below the sea surface and, based on the atmospheric pressure acting on the sea surface (P_{Aw}), the mean density in the water column (ρ_w) and the gravitational acceleration (g). This pressure can be converted to sea level height using the relationship between measured pressure (P) and depth (D) as given by, e.g. Pugh (1987):

$$P = P_{Aw} + \rho_w g D \quad (5)$$

As an example of a subsurface pressure tide gauge system, Figure 6 (reproduced from Pugh (1987)) shows a schematic drawing of a so-called bubbler gauge. The cylinder is open at the bottom so that water can flow in. A steady flow of compressed air or other gas is let into the connecting tube and can bubble out through an orifice (the pressure-point). As explained by Pugh (1987), for low rates of gas escape, gas pressure equals water pressure (P). Apart from some small pressure gradients in the connection tube, this pressure is transmitted along the tube and recorded by the recording system, see Pugh (1987) for a description of recording systems.

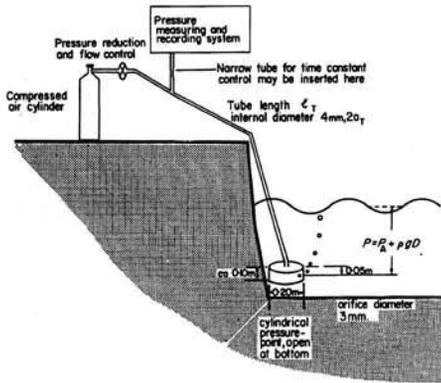


Fig. 6. Subsurface pressure tide gauge system; reproduced from Pugh (1987)

Equation 5 requires atmospheric pressure measured at the instantaneous sea surface. Since this pressure cannot be obtained in practice, atmospheric pressure is measured by a sensor at a height (h_a) above the sea surface. According to Carrera *et al.* (1996) the relation between atmospheric pressure (P_A) at height h_a and atmospheric pressure at the instantaneous sea surface (P_{A_w}) is

$$P_{A_w} = P_A g h_a \rho_a \quad (6)$$

in which ρ_a is the density of the air between the sensor and the sea surface.

From equation 5 and equation 6 it is clear that the precision of sea level heights determined by subsurface pressure tide gauges depend on the precision of the measurements of atmospheric pressure and "water" pressure, the precision of the value used for the gravitational acceleration, the precision of the values used for resp. the water and air density, and on the precision of the determined height between the pressure point in the cylinder and the atmospheric pressure sensor.

Carrera *et al.* (1996) give a comprehensive description of these error sources contributing to the overall precision of determined sea level heights. Many of these error sources themselves depend on measurements of a number of other parameters, e.g., water density is a function of salinity, water temperature and pressure. According to Carrera *et al.* (1996), of the six parameters involved, uncertainties in "water" pressure, water density and gravity contribute the most to the resulting uncertainty in sea level height, whereas the contribution of uncertainties in atmospheric pressure, atmospheric density and height between the pressure point in the cylinder and the atmospheric pressure sensor are relatively small.

Pugh (1987) estimates that the pneumatic system in a subsurface pressure tide gauge with a connecting tubes up to 200 meters, can produce water head equivalents (i.e., the depth of water which would produce a specific pressure) with a precision of 1 cm.

Besides random errors, the accuracy of determined sea level heights is also influenced by the occurrence of systematic errors. For example, as described by Xu (1990), systematic errors are introduced by the non-linear relationship between water depth and pressure; hydrostatic pressure ranges between 0.1019 and 0.1034 kg/cm². Furthermore, large errors occur when water is forced into the connection tube by waves, for more details see Pugh (1987).

One of the major advantages of subsurface pressure tide gauges is that they are relatively convenient to use and they can be operated under difficult environmental conditions. One of their major drawbacks is that it is often difficult to relate the zero-height point of the system to the land based tide gauge bench mark. This is due to (different) biases and drifts inherent to the air pressure sensor and the "water" pressure sensor. This problem is addressed in detail by Woodworth *et al.* (1996), who estimate that subsurface pressure tide gauge data is often related to the land datum with a precision of about 2 cm. They mention the following methods which are presently used to overcome this datum problem:

- simultaneous measurements at a nearby stilling well
Although this method works fairly well as long as comparisons are performed based on several complete tidal cycles to remove the effect of any lag in the well, as described in Section 2.3, stilling wells will introduce systematic errors of their own.
- measurements of tide poles or tide poles with float by an observer
Read outs can only be made at a limited number of times and accuracy of measurements depends, a.o., on the state of the sea surface, see Sections 2.1 and 2.2
- water level "switches" in mini-stilling wells
Although these switches show great promise, presently they are not able to entirely eliminate the effect of waves, and are probably accurate to only a few cm
- using "comparators", or precisely calibrated reference pressure devices
Although they appear to give a precision of at least 1 cm for the datum control, they do not provide near-continuous datum check, are clumsy to operate and are not well documented.

All the above is based on Woodworth *et al.* (1996).

Woodworth *et al.* (1996) describe a method which seems to be able to provide datum control with a precision in the order of only 1 mm. This method is based on an additional pressure point situated at a known height approximately near mean sea level. A description of this method is beyond the scope of this article and interested readers are referred to Woodworth *et al.* (1996).

2.6 Open sea pressure gauges

The tide gauge instruments as described in the preceding sections can only be used to measure sea level variations relative to a tide gauge bench mark on land, or on an off-shore platform. Using so-called open sea pressure gauges it is also possible to measure sea level heights at sites far away from coasts. Usually, these tide gauge systems consist of a pressure sensor, placed on the ocean bottom, which is able to measure and record the pressure of the overlying water column. Depending on the water depth in which the tide gauge has to operate, different systems have been developed. In principle there are two categories of open sea pressure gauges: deep sea pressure gauges and pressure gauges for use on continental shelves.

According to Pugh (1987), in relative shallow water the tide gauge is often attached to a ground line, which simplifies the recovering of the gauge. In shallow water it is also possible to transmit recorded data to, e.g., a surface buoy, either through a cable or acoustically. This has the major advantage that the sea level data can be collected without removing the tide gauge from its location on the bottom of the sea. Consequently, data can be available almost real-time and tide gauges can continue operating on the same location (at least for as long as its power supply lasts).

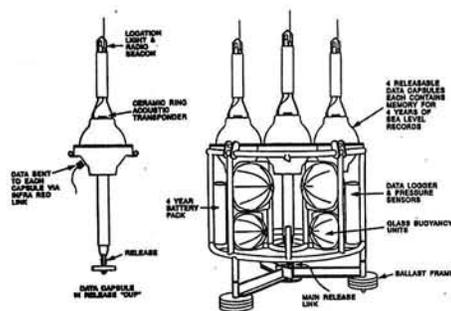


Fig. 7. Deep sea pressure tide gauge system capable of 4 years of recording; reproduced from Smithson (1997)

For deep sea operations, pressure tide gauges are usually built into a protective framework which is lowered to the bottom of the sea. These pressure gauges can operate in depths of over 4000 meters, with a resolution of 0.01 meters; see Pugh (1987). Since, presently, no methods have been developed to transmit the data while the pressure tide gauge is situated at the bottom of the sea, after a certain amount of time (usually 1 year), the tide gauge has to be recovered. This has the disadvantage that sea level data is only available after recovery of the tide gauge, and undisturbed time series are only available over relative short periods of time. After recovery, the tide gauge can be lowered back to the bottom of the sea, but it will always be in a (slightly) different location and the measurements series is discontinued for the time period needed for the recovery and replacement of the tide gauge.

A new deep sea pressure tide gauge called "MYRTLE" (Multi Year Return Tide Level Equipment) can, partly, overcome these problems. Figure 7 shows a schematic drawing of this tide gauge system which consists of four data logger capsules. After one year of operation the first capsule is released and floats to the ocean surface where it can be recovered. After two years of recording the second capsule is released, etc. To prevent loss of data, each capsule contains data for the full measurement series as performed up to the release of the capsule. After four years of operation, the whole framework is released and can be recovered.

Major limitations for the accuracy of open sea pressure system are the accuracy of the pressure sensor itself, drift of the pressure transducer zero, and settlement or lift of the instrument relative to the bottom of the sea. Pugh (1987) estimates that the framework will settle into the sediments with a velocity of a few cm/month. The accuracy of the pressure transducer is, a.o., affected by its sensitivity to temperature changes. According to Pugh (1987), in depths of (at most) 200 meters drift of the transducer zero can be reduced to a few cm/month, while in depths of around 4000 meters drift can be significantly larger. In addition, according to Banaszek (1985), the pressure sensor, and especially the framework in which the tide gauge is mounted, can seriously distort the velocity field and, consequently, pressure detected by the sensor deviates from the hydrostatic pressure which is related to depth. Errors in observed pressure ranging between 1 and 30 mb have been found; see, e.g., Banaszek (1985) or Muir (1978).

Variations in density through the water column above the pressure transducer should be taken into account as well, e.g., by calibration of the instrument based on measurements through the water column. Furthermore, pressure as measured by the transducer not only depends on the

amount of water overlying the transducer but also on the atmospheric pressure acting on the sea surface. Rae (1976) estimates that pressure variations due to atmospheric pressure can be in the order of 50 millibars. However, since sea level response to variations in atmospheric pressure is almost inverse barometric (at least for "open" seas, i.e., away from continental boundaries), the total pressure at the sea bottom will not be significantly affected. On the other hand, according to Rae (1976), density variations (e.g., due to internal waves and the formation of thermoclines) can cause variations in total bottom pressure of about a few millimeters.

3 Sampling rate and averaging method of tide gauge readings

In the preceding section, for the six major tide gauge systems, error characteristics have been described. However, predictions of sea level variation curves are usually based on some kind of average values (e.g., monthly or yearly sea level heights) instead of on instantaneous measurements of sea level height. Consequently, the accuracy of these predictions is not only determined by the accuracy of the individual tide gauge measurements but also by the sampling rate of these measurements and the method applied to these individual measurements to form mean values.

The reason for basing evaluations of sea level variations on mean sea level heights is simply that records of instantaneous sea level values are, in general, not available. In the past, when observations of sea level height were simply written down, only monthly mean values (or even yearly mean values) were recorded for long-term keeping. Even after the introduction of mechanical recording (e.g., on a paper chart mounted on a drum) usually only monthly mean values were stored. For the more recent past (last few decades ?) sometimes hourly mean values are available, often obtained by digitizing the paper charts at hourly intervals. Although state-of-the-art tide gauge systems usually work with a relatively high sampling rate, again, only average values are stored. As an example, in the case of a reflection tide gauge (see Section 2.4) measurements are made at a rate of 1/s but only 3-minute average values are stored.

An advantage of mean sea level heights is that, depending on the method used to form the averages and the time span over which the average is taken, high frequency signals (like waves) and periodic effects (like tides) are, partly, removed from the data. In the following sections different averaging methods will be described which have been used through the years. In complexity these methods range from simply taking the arithmetic mean of only two measurements a day, to low-pass filtering based on 24 hourly values a day.

A limitation of historic sea level data based on manual readings by an observer is that these data often have a relatively low sampling rate, i.e., mean values are based on only a few readings a day. In addition, the oldest data is often based on irregular sampling over the tidal cycle, e.g., a tide pole was only read during daytime (at high and low water). Since this limited number of sea level height samples outlines the averaging method which can be used, its effects will be described in the following sections where the different averaging methods will be introduced.

Finally, it should be noted that the accuracy of mean sea level values is also influenced by how well outliers and systematic errors in the individual sea level measurements have been corrected for. As described in detail by Pugh (1987), prior to forming the average values, individual measurements should be checked for reading errors and corrected for scaling errors (due to e.g., timing errors in the control clock) and it might be necessary to use interpolated values to fill gaps in the data series. Since this is beyond the scope of this article, for a description of methods which can be used to check the recorded measurements, the interested reader is referred to, e.g.,

Pugh (1987).

Depending on the averaging method used, periodic effects (with a period which is short relative to the time span over which the measurements are averaged) are, partly, removed from the resulting mean values. However, the influences of low-frequency tides (like the lunar nodal tide with a period of 18.6 year) will more or less remain in monthly and even in yearly mean sea level heights. In addition, mean monthly and annual sea levels are influenced by the occurrence of storm surges.

3.1 Low-pass filtering of hourly values

The most advanced method of forming mean (monthly or annual) sea level heights is low-pass filtering of hourly (mean) sea level heights to obtain smoothed daily mean sea levels. By taking the arithmetic mean of these smoothed daily sea levels over a period of a month or a year, resp. monthly or annual mean sea level heights can be derived.

Different low-pass filters, requiring a different number of hourly (mean) sea level values have been developed, the most widely used is the so-called Doodson X_0 filter which uses 39 hourly values. As described in, e.g., IOC (1985) the smoothed daily value ($X_F(t)$) is derived by applying the filter F_m to the 39 hourly values in the symmetric window around time t :

$$X_F(t) = \frac{1}{30} F_0 X(t) + \frac{1}{30} \sum_{m=1}^{19} F_m [X(t+m) + X(t-m)] \quad (7)$$

in which the filter elements are defined by

m	0	1	2	3	4	5	6	7	8	9	10	11	12	13	14	15	16	17	18	19
F_m	0	2	1	1	2	0	1	1	0	2	0	1	1	0	1	0	0	1	0	1

This Doodson low-pass filter, and other filters based on even more hourly observations (e.g., 72, or even 169 values), are designed to yield an optimal elimination of the influence of diurnal and semi-diurnal tidal constituents on the resulting monthly and annual mean sea level heights.

3.2 Arithmetic mean of hourly values

Instead of using low-pass filters to smooth the hourly values in a first step, often the hourly values themselves are used directly to form monthly or annual mean sea level heights. These mean values are simply determined as the arithmetic mean over all (mean) hourly sea level values in resp. a month or year. According to Xu (1990) this method is the most widely used because it involves less computational effort and, since the major part of the effects of the diurnal and semi-diurnal tides are removed, the resulting mean values are almost similar to those produced by low-pass filtering.

From 1971 onwards in the Netherlands mean (monthly and annual) sea level heights have been determined as arithmetic means of hourly sea level heights; see Van der Hoek Ostende and Van Malde (1989). At first these hourly values were derived by digitizing mechanical produced sea level charts at hourly intervals. Starting in 1987, mechanical recording on charts has been replaced by digital recording of 10-minute mean values. The resulting hourly values are rounded off to cm-level, the arithmetic means over all hourly values in a year are rounded to the mm-level.

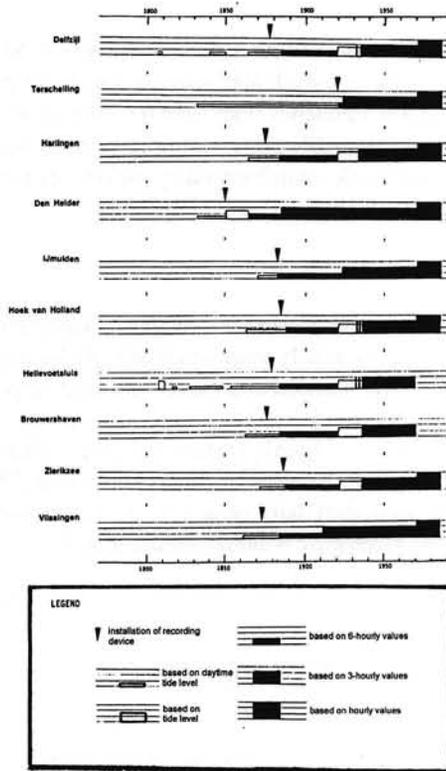


Fig. 8. Historical overview of averaging methods applied to tide gauge measurements to form mean annual sea level height; reproduced from Van der Hoek Ostende and Van Malde (1989)

Figure 8, reproduced from Van der Hoek Ostende and Van Malde (1989), shows for 10 tide gauges in the Netherlands for which more than 100 year of tide gauge data is available, which method of forming annual mean sea level heights has been applied through the years. This Figure clearly shows that, the further back in time, the smaller the number of daily observations used to determine annual mean sea level heights. In the following sections, in descending order of number of daily observations, these historic methods of forming annual mean sea levels will be described.

3.3 Arithmetic mean of 3-hourly values

Prior to using hourly measurements, annual mean sea level heights were based on 8 measurements/day. As described by Van der Hoek Ostende and Van Malde (1989), sea level values at 2, 5, 8, 11, 14, 17, 20, and 23 o'clock were read from the paper charts and rounded off at cm-level. From these 3-hourly values annual mean sea level was simply determined as the arithmetic mean over all values in one year, rounded to 1 mm.

As described by Van der Hoek Ostende and Van Malde (1989) measurements performed before 1961 usually refer to Amsterdam time, later measurements are usually made with reference to MET (Middle European Time), the time difference between the two systems being 40

minutes. However, according to Van der Hoek Ostende and Van Malde (1989) the effect of this time-shift on determined annual mean values is negligible.

Prior to using 3-hourly measurements, annual mean sea level heights were (usually) based on only 4 daily measurements, see Section 3.4. As can be seen from Figure 8, the moment when the change was made between these two sampling rates varies widely for the various stations. As an example, for the station Den Helder already around 1885 mean sea levels were based on 8 observations/day, whereas only from 1936 onwards this method was used for all tide gauge stations along the Dutch coast.

Experiments based on hourly mean values for the period between 1971 and 1986, as performed by Van der Hoek Ostende and Van Malde (1989), showed that the difference between annual mean sea level heights based on the 3-hourly values and annual mean sea level heights based on hourly values is relatively small. The mean value for the difference between these two time series of annual values ranges between -0.4 and +0.6 mm for the various stations, whereas the standard deviations range between 0.3 and 1.0 mm. According to Van der Hoek Ostende and Van Malde (1989), the stations under consideration are situated in largely varying tidal regimes, which can be assumed representative for the Dutch coastal zone. Therefore, they conclude that using 3-hourly values instead of hourly values will hardly introduce errors for other tide gauge stations as well.

3.4 Arithmetic mean of 6-hourly values

From Figure 8 it can be seen that after mechanical recording became available, for most stations the switch was made from annual mean sea levels based on (daytime) tide levels to annual mean sea level heights based on 6-hourly values. Analogous to the method as applied to the 3-hourly measurements, sea level heights at 2:00, 8:00, 14:00, and 20:00 o'clock Amsterdam time were read from the paper chart and the arithmetic mean of all values in one year (rounded at the mm-level) was taken to form annual mean sea level heights.

Van der Hoek Ostende and Van Malde (1989) have compared, for the period between 1971 and 1986, annual mean sea levels based on 6-hourly values with those based on hourly values. For a number of stations the difference between the two resulting time series is rather large. These are all stations situated in a rather special tidal regime, where a short-duration increase in sea level height is seen at low tide, which is related to a relative long duration of low tide. Consequently, relative large deviations are introduced for these stations if annual mean sea level is based on 6-hourly values instead of on hourly values.

Also for stations situated in more "normal" tidal regimes there are significant differences between annual mean values based on 6-hourly values and those based on hourly values. For these stations, the mean value for the difference between the two time series of annual values range between -3 and +6 mm, whereas the standard deviations range between 1.2 and 2.8 mm; see Van der Hoek Ostende and Van Malde (1989).

3.5 Mean sea level heights determined from mean tide level

For a number of stations, between 1920 and 1935, yearly mean sea level heights were based on mean tide levels; see Figure 8. Since hourly values were, in principle, available this was probably done in order to simplify computations. Contrary to the period before the installation

of mechanical recording devices (see next section), these tidal mean values were not only based on daytime measurements of high and low tide level but on nighttime values as well. These high and low tide values were simply read from the paper charts with a resolution of 1 cm; see Van der Hoek Ostende and Van Malde (1989) for more details. Until 1971, in addition to hourly sea level values, these high and low water values have been recorded.

By taking the arithmetic mean of resp. all high tide values and all low tide values in one year, annual mean high tide and low tide values were determined and rounded to the nearest mm. By averaging these two mean tide values the annual mean tide value was found. From these annual mean tide levels, annual mean sea levels were determined by applying a correction factor. According to Van der Hoek Ostende and Van Malde (1989) this was based on the, false, assumption that with a high enough accuracy the difference between annual mean sea level and annual mean tide level for a specific station could be assumed constant.

The corrections as applied to the mean tide level were different for every station, for the stations under consideration they range between -15.5 and +18 cm. Unfortunately, although the values used for these corrections are known, it is (as yet) unknown how these values have been derived, see Van der Hoek Ostende and Van Malde (1989) for more details.

For the period between 1971 and 1986, Van der Hoek Ostende and Van Malde (1989) compared annual mean sea levels based on mean tide values with annual mean sea levels based on hourly values of sea level height. For the various stations they found a wide range in differences between these two time series of annual mean values. The mean value for the differences ranging between -25 and +21 cm, and the standard deviations ranging between 0.4 and 1.5 cm.

In addition, experiments by a number of authors have shown that the tidal regime along a major part of the Dutch and German coasts is changing, i.e., the mean difference between high and low tide increases. Consequently, estimating mean sea level height by applying a, constant, correction to mean tide level can yield an (increasingly) large systematic error. As a result of this phenomenon and the large standard deviations found between time series based on mean tide values and hourly sea level values, Van der Hoek Ostende and Van Malde (1989) conclude that annual mean sea level values based on mean tide values are rather inaccurate.

3.6 Mean sea level heights determined from mean daytime tide level

As described by Van der Hoek Ostende and Van Malde (1989), in the Netherlands, before the installation of mechanical recording devices only the high and low water level were measured. Measurements were made during daytime (between 6:00 and 18:00 Amsterdam time), with a 1-cm resolution. All values of high water level in one year were combined in a mean daytime high-tide level. Mean daytime low-tide level was accordingly determined as the arithmetic mean of all daytime low tide values in a year. Both mean tidal values were determined with a resolution of 1-cm. Next, annual mean daytime tide level was calculated as the arithmetic mean of mean daytime high-tide level and mean daytime low-tide level, and rounded off to an integer number of cm. This rounding off was performed as follows, if the integer part of the mean value was an even number it was rounded down, an uneven value was rounded up.

From these annual mean daytime tide levels the annual mean sea levels were determined by applying a correction. The various stations have different correction factors which are, usually, constant over the whole period of time for which this method was applied. However, for some stations the value used for the correction factor has been changed at a certain time. For the 10 stations under consideration, values for the correction factor range between -17 and +18 cm.

It should be noted when comparing sea level data for various tide gauges, that although for most tide gauges (as mentioned in these sections) the switch to 6-hourly measurements was made around the same time (somewhere around 1885) for station Terschelling annual mean sea levels were still based on annual mean daytime tide level until around 1920, see Figure 8.

For the period between 1923 (resp. 1936) and 1960, Van der Hoek Ostende and Van Malde (1989) compared, for 7 stations along the Dutch coast, annual mean sea level heights based on mean annual daytime tide level with those based on mean annual tide level. They found differences between the mean values of these two time series of up to 6.4 mm (for station Den Helder). However, as explained by Van der Hoek Ostende and Van Malde (1989), this difference in mean value does not necessarily have a large influence on the accuracy of the mean daytime tide level method. Since both time series of annual mean sea level are obtained by applying a correction to the mean tide levels, the difference between the two series can be minimized by changing the correction factor as applied to the mean annual daytime tide values. However, standard deviations for the time series based on mean daytime tide level are between 10 and 30% higher than standard deviations for the series based on mean tide level; both relative to annual mean sea level heights based on 3-hourly measurements.

As described by Van der Hoek Ostende and Van Malde (1989) in addition to inaccuracies introduced by using only daytime samples, less precise annual mean sea levels are obtained because annual mean daytime high tide and low tide levels are rounded off to an integer number of centimeters. Compared to annual mean sea level heights based on mean high tide and low tide values, which are rounded to the nearest mm, this rounding off already introduces a standard deviation of almost 3 mm.

4 Conclusions and recommendations

In the preceding sections error characteristics of a number of tide gauge systems have been discussed. It is clear that tide gauges without automatic recording are unsuitable for high quality monitoring of sea level variations over a longer period of time, due to the limited measuring frequency and the susceptibility to reading and recording errors. State-of-the-art tide gauge systems are equipped with a mechanical or digital recording device which allows continuous measurements or measurements with a relative high sampling rate. These recording devices introduce some errors of their own, e.g., errors due to mechanical problems (like friction) or due to aliasing. In addition, the quality of the measured sea level heights is influenced by the occurrence of systematic errors. Tide gauges based on different techniques, i.e., based on a float in a stilling well, acoustic travel times, or a comparison between water pressure and atmospheric pressure, are susceptible to different systematic errors.

For a new site, a suitable tide gauge system should be selected based on characteristics of the specific location (e.g., is it easy to reach for maintenance), environmental conditions (e.g., the occurrence of large currents), expected measuring precision, and available budget. However, for analysing existing tide gauge data, it should be remembered that often less optimal tide gauge systems have been used. For example, in the Netherlands tide gauge data from the 19th century is usually based on manual recording.

Since often for longer periods of time only monthly or annual sea level heights are available, the method used to form these average values should be considered as well. Especially for relatively old tide gauge data, averaging methods have been used which yield biased mean values. For

example, in the Netherlands only from 1935 onwards for all tide gauges mean values were based on 3-hourly or hourly values. Prior to this date, often only 6-hourly values were used or mean sea level heights were even based on mean tide levels. When evaluating long series of sea level values, differences in tide gauge systems and averaging methods used, occurring over the time span of the measurements should be taken into account.

References

- Banaszek, A. (1985). Procedures and problems associated with the calibration and use of pressure sensors for sea level measurements. In *Evaluation, comparison and calibration of oceanographic instruments*, volume 4 of *Advances in underwater technology and offshore engineering*. Graham and Trotman, London.
- Carrera, G., Tessier, B., and O'Reilly, C. (1996). Statistical behavior of digital pressure water level gauges. *Marine Geodesy*, **19**, 137–163.
- Christensen, E., Haines, B., Keihm, S., Morris, C., Norman, R., Purcell, G., Williams, B., Wilson, B., Born, G., Parke, M., Gill, S., Shum, C., Tapley, B., and Nerem, R. S. (1994). Calibration of TOPEX/POSEIDON at platform Harvest. *Journal of Geophysical Research*, **99**(C12), 24465–24485.
- Cross, R. (1968). Tide gauge frequency response. *Journal of waterways and harbour division*, **94**(WW3). American Society of Civil Engineers.
- Diamante, J., Pyle, T., Carter, W., and Scherer, W. (1987). Global change and the measurement of absolute sea-level. *Progress in Oceanography*, **18**, 1–21.
- Hamon, B. and Godfrey, J. (1980). Mean sea level and its interpretation. *Marine Geodesy*, **4**(4), 315–329.
- IOC (1985). Manual on sea level measurement and interpretation. Technical Report 14, Intergovernmental Oceanographic Commission.
- Lennon, G. (1970). Sea level instrumentation, its limitations and the optimization of the performance of conventional gauges in Great Britain. In R. Sigl, editor, *Coastal Geodesy; symposium, Munich, July 1970*. z. uitg. International Union of Geodesy and Geophysics; Technical University Munich, Institute for Astronomical and Physical Geodesy.
- Martin, D., Chapin, J., and Maul, G. (1996). State-of-the-art sea level monitoring. *Marine Geodesy*, **19**, 105–114.
- Montag, H. (1970). On the accuracy of determination of secular variations of mean sea level at the Baltic Sea coast. In R. Sigl, editor, *Coastal Geodesy; symposium, Munich, July 1970*. z. uitg. International Union of Geodesy and Geophysics; Technical University Munich, Institute for Astronomical and Physical Geodesy.
- Muir, L. (1978). Bernoulli effects on pressure-activated water level gauges. *International Hydrographic Review*, **55**(2), 111–119.
- Noye, B. (1974). Tide-well systems: 3. Improved interpretation of tide-well records. *Journal of Marine Research*, **32**, 183–194.
- Pugh, D. (1987). *Tides, surges and mean sea level, a handbook for engineers and scientists*. John Wiley & sons.
- Rae, J. (1976). The design of instrumentation for the measurement of tides offshore. In *The Hydrographic Society Symposium on tide recording, proceedings*, number 4 in special publication Hydrographic Society. Hydrographic Society.
- Sager, G. and Matthäus, N. (1970). Theoretical and experimental investigations into the damping properties of tide gauges. In R. Sigl, editor, *Coastal Geodesy; symposium, Munich, July 1970*. z. uitg. International Union of Geodesy and Geophysics; Technical University Munich, Institute for Astronomical and Physical Geodesy.
- Smithson, M. (1997). Tide gauges. presented at the summerschool: sea level changes on micro to macro time scales: measurements, modelling, interpretation and application; Kos, Greece.
- Van der Hoek Ostende, E. and Van Malde, J. (1989). De invloed van de bepalingwijze op de berekende gemiddelde zee­stand. Technical report, Ministerie van verkeer en waterstaat, dienst getijdewateren. nota GWAO-89.006, in Dutch.
- Van der Made, J. (1987). *Analysis of some criteria for design and operation of surface water gauging networks*. Ph.D. thesis, Delft University of Technology. Van Gorcum, Assen.
- Woodworth, P., Vassie, J., Spencer, R., and Smith, D. (1996). Precise datum control for pressure tide gauges. *Marine Geodesy*, **19**, 1–20.
- Xu, P. (1990). Monitoring sea level rise. Technical report, Faculty of Geodetic engineering, TU Delft.

Research plan and progress report for the TOPEX/POSEIDON extended mission

Marc Naeije, Roger Haagmans, Ernst Schrama, Karel Wakker, and Remko Scharroo

Abstract

This paper covers the highlights of the research plan for the TOPEX/POSEIDON Extended Mission (TPEM). It partly elaborates on themes of the original T/P science investigation plan, and partly deals with new challenges for ocean circulation, height systems, and gravity field studies. In addition, some recent results and publications are summarized.

1 Introduction

In 1996 the Announcement of Opportunity (AO) for the TOPEX/POSEIDON Extended Mission (TPEM) was issued as a follow-on to the TOPEX/Poseidon (T/P) Science Investigation plan. It embroidered on the successes of the exploitation of altimeter data from the T/P mission for studying ocean circulation and ocean tides. DEOS, the Delft Institute for Earth-Oriented Space Research, a joint venture between the faculties of Aerospace and Geodetic Engineering of the Delft University of Technology, has a long record in altimetry related studies. It took advantage of the AO by updating its research plan. This plan partly elaborates on themes of the original T/P science investigation plan, and partly deals with new challenges.

In principal, the plan is based on three lines of interest: the stationary behavior of the global mean sea surface, temporal variations of the sea surface at various scales, and calibration and validation of altimeter data. DEOS is most interested in continuing to use altimetry in its long term research activities. The justification is a continuation of an important measurement series of the global oceans. DEOS strives to use this series for continuing to understand the long periodic behavior of sea level anomalies which are vital for observing and predicting ENSO events and inferred Rossby waves, annual and semi-annual cycles in the oceans and global ocean tides. Obviously, continuation of the T/P time series implies improved results for ocean tide models, mean sea surface, ocean circulation, and sea level change. The activities in the field of orbit determination are drastically reduced due to the present accuracy of T/P orbits. The main focus is on exploring new possibilities of altimeter and TDRSS case studies for orbit and gravity field improvement, and on orbit improvement of LEO satellites from T/P ocean tide models.

Other regions of interest are the European continental shelf, the Indonesian Archipelago, and the Chinese Sea. Here T/P results are used for regional ocean tide and circulation models,

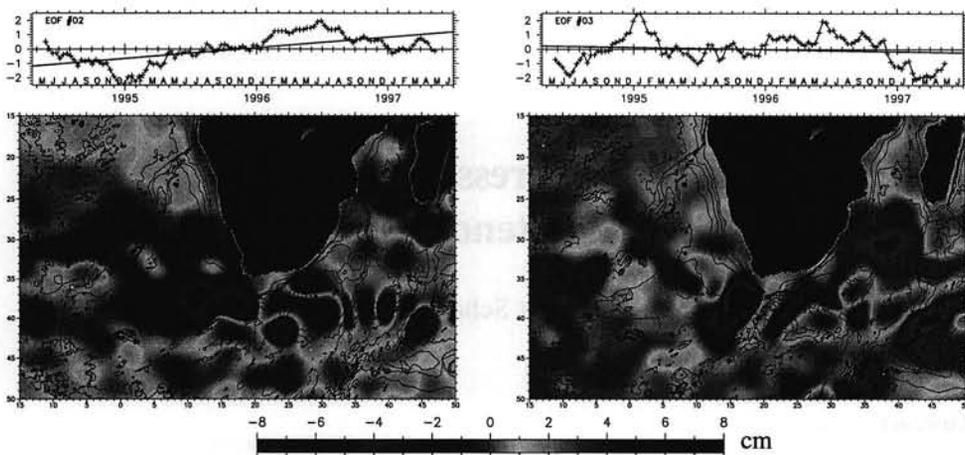


Fig. 1. 2nd and 3rd EOF of T/P altimeter data (1994-1997) for the oceans around South Africa.

connection and unification of height datums, regional sea level change and land subsidence studies.

The TPDM project at issue unites three Dutch institutes and one German institute, viz. the Delft University of Technology (PI: K. F. Wakker, CoIs: M. C. Naeije, R. Scharroo, R. H. N. Haagsmans, and E. J. O. Schrama), Utrecht University (CoIs: W. P. M. de Ruijter, and P. J. van Leeuwen), Survey department Rijkswaterstaat (CoI: R. C. V. Feron), and the Technical University of Munich (CoI: R. Rummel).

2 Ocean circulation

The extension of the high precision altimetric database has been used to obtain better ocean tide, and ocean circulation results, both global and regional. In particular T/P altimetry is well suited to regional tide model improvement and storm surge predictions for the North Sea and the Chinese Sea. Intra- and inter annual variations, and the eddy shedding mechanisms of the western boundary currents have been studied in detail. A long term objective is to determine the decadal variability in the climate system. A key element in this system is the global thermohaline circulation in which the Agulhas region is thought to be a major link. For a better understanding of the Agulhas system the precise mechanism of shedding of rings from the Agulhas Current is studied by means of assimilation of altimetry data in a regional ocean model. Also the decay and the interaction of Agulhas rings with the bottom topography are studied using a combination of altimeter data, infrared data, and optical data together with numerical models. To illustrate the complexity of the dynamics in this area, the 2nd and 3rd EOF of T/P data (1994-1997) for the oceans around South Africa are plotted in Figure 1, which reveals a mix of short (3 months) to long periodic (3 years) patterns.

In addition, assimilation studies are also applied to the equatorial Pacific region to obtain a better insight in the El Niño/Southern Oscillation, and ocean-atmosphere interaction. Special attention is given to the impact of the upper ocean salinity structure on the western Pacific dynamics. For this purpose the ability of the altimeter to serve as a salinometer is studied. This

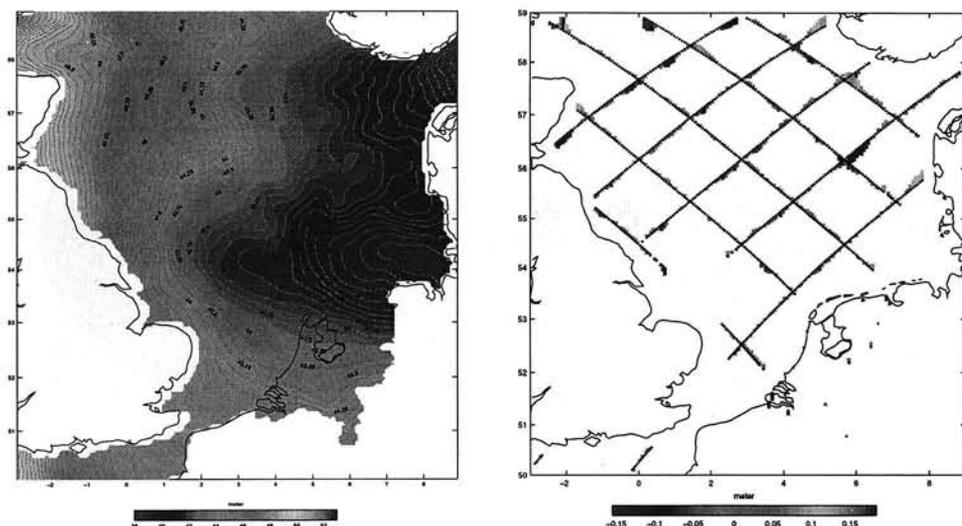


Fig. 2. *left:* North Sea geoid GEONZ97 w.r.t. GRS80 in meters (contour interval 0.25 cm). *right:* Differences of T/P and NEREF (GPS/leveling) data points with North Sea geoid; $\mu = -0.1\text{cm}$, $\sigma = 4.2\text{cm}$

research is carried out in cooperation with the Royal Dutch Meteorological Institute (KNMI), NOAA, and NCEP.

3 Height systems

On a global scale the accuracy of sea level change estimates is improved. On regional scales the combination of T/P with ERS data and/or in situ data may improve the separability between the T/P instrument drift and sea level change. The current accuracy of T/P data enables incorporation of altimetry in the North Sea Sea level monitoring system (NOSS) for sea level rise and land subsidence studies.

The precise determination of the sea surface in the region of the Indonesian Archipelago enables the connection and unification of height datums of the Indonesian islands, and analysis of the steady-state ocean topography. This is an important step towards an ultra precise regional sea-land monitoring system based upon various terrestrial and satellite observation techniques comparable to NOSS. Figure 2 illustrates a successful attempt to connect a land-sea based gravimetric geoid for the North Sea with the geoid at mean sea level by fitting a correction model through T/P data at sea and GPS and leveling data on land.

4 Gravity field

Another aim is the determination of a highly detailed precise mean sea surface using T/P data as a reference for ERS, and Geosat ERM and GM data. Then from this mean sea surface an improved geoid can be obtained whenever detailed models or data become available of the mean ocean dynamic topography. Furthermore, gravity anomalies, and gravity gradients are

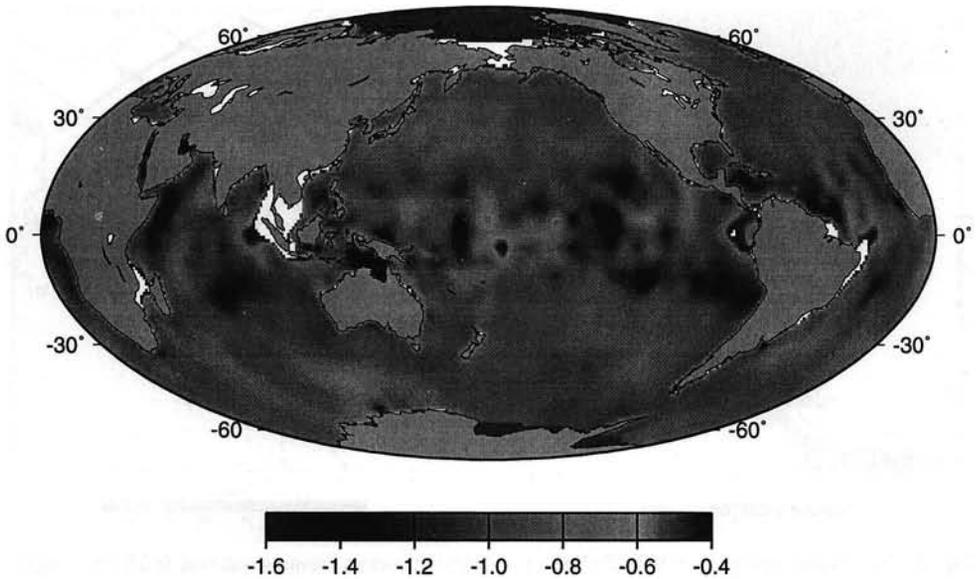


Fig. 3. IB regression coefficient in cm/mbar derived from T/P collinear differences (cycles 2-84).

determined for regional geophysical interpretation, and improvement of statistical models of the global gravity field.

5 Outreach

In behalf of national and international projects DEOS is working on a consistent altimeter data base from all past and present altimeter missions, taking into account all different references and corrections. Validation of these corrections is a key issue. Figure 3 illustrates that the regression between pressure field and T/P height anomaly is geographically dependent and on average not equals -1 cm/mbar. At least, one has to take the variations in wind speed and stress into account.

As pointed out earlier DEOS is interested in continuing its altimetry related research. This means that already preparations are being made for incorporating the data from the forthcoming Jason mission. A proper phasing of the orbit of Jason w.r.t. T/P can lead to an enhanced sampling strategy for mapping meso-scale variability phenomena and ocean tides. Intersatellite calibration with respect to T/P and ERS-2 or other orbiting altimeters will play a significant role in these studies.

References

- Bruijne, A. J. T. de, R. H. N. Haagmans, and E. J. de Min (1997). *A preliminary North Sea Geoid model GEONZ97*. MD report MDGAP-9735, Survey Department, Rijkswaterstaat, Netherlands.
- Feron R. C. V., W. P. M. de Ruijter, and P-J. van Leeuwen (1997). A new method to improve the mean sea surface topography from altimeter observations. Accepted for publication in *J. Geophys. Res.*

- Haasbroek N. (1996). *Characteristic of the Gulf Stream from TOPEX/Poseidon*. MSc. thesis, Delft University of Technology, Faculty of Geodetic Engineering, 86 pp.
- Khafid (1997). *On the unification of Indonesian local height systems*. Deutsche Geodätische Kommission, Reihe C, Heft 488
- Naeije M. C. and R. Scharroo (1996). *A global sea level variability study from almost a decade of altimetry*. In: *Digest International Geoscience and Remote Sensing Symposium (IGARSS)*, Volume I, IEEE Catalog Number 96CH35875, Lincoln, Nebraska.
- Naeije M. C., R. Scharroo, and K. F. Wakker (1996). *Monitoring of Ocean Current Variations*. NUSP project 1.2/OP-04, NUSP report 96-08 (ISBN 90-5411-190-9), Netherlands Remote Sensing Board (BCRS), Delft, The Netherlands.
- Naeije M. C. and K. F. Wakker (1997). *Global Analysis of Sea Surface Height and Temperature*. *Proceedings of the Third ERS Symposium*, ESA SP-414, Florence (in press).
- Ruijter, W. P. M. de, P.-J. van Leeuwen, and J. R. E. Lutjeharms (1997). *Generation and evolution of Natal Pulses: solitary meanders in the Agulhas Current*. submitted for publication in *J. Phys. Ocean.*
- Scharroo R. and P. N. A. M. Visser (1997). *Precise orbit determination and gravity field improvement for the ERS satellites*. Accepted for publication in *J. Geophys. Res.*
- Schrama E. J. O. (1997). *A study of the global oceanic response to pressure and wind variations by application of satellite altimetry*. *EOS*, G22A-6, April 29.

Faint, illegible text, possibly bleed-through from the reverse side of the page. The text is too light to transcribe accurately.

Progress Letters of the Delft Institute for Earth-Oriented Space Research:

97.1

98.1

1000 University Avenue, Chicago, Illinois 60607-7073

Phone: (773) 936-3200 Fax: (773) 936-3201

Internet: <http://www.library.uchicago.edu>

Library Hours: Monday - Friday, 9:00am - 5:00pm

Saturday, 10:00am - 4:00pm

Sunday, 12:00pm - 4:00pm

Special Hours: See www.library.uchicago.edu

For more information, contact your librarian

or call (773) 936-3200

or visit www.library.uchicago.edu

Library of The University of Chicago

1000 University Avenue, Chicago, Illinois 60607-7073

Phone: (773) 936-3200 Fax: (773) 936-3201

Internet: <http://www.library.uchicago.edu>

Library Hours: Monday - Friday, 9:00am - 5:00pm

Saturday, 10:00am - 4:00pm

Sunday, 12:00pm - 4:00pm

Special Hours: See www.library.uchicago.edu

For more information, contact your librarian

or call (773) 936-3200

or visit www.library.uchicago.edu

Library of The University of Chicago

1000 University Avenue, Chicago, Illinois 60607-7073

Phone: (773) 936-3200 Fax: (773) 936-3201

Internet: <http://www.library.uchicago.edu>

Library Hours: Monday - Friday, 9:00am - 5:00pm

Saturday, 10:00am - 4:00pm

Sunday, 12:00pm - 4:00pm

Special Hours: See www.library.uchicago.edu

For more information, contact your librarian

or call (773) 936-3200

or visit www.library.uchicago.edu

Library of The University of Chicago

1000 University Avenue, Chicago, Illinois 60607-7073

Phone: (773) 936-3200 Fax: (773) 936-3201

