

Least-Squares Estimation of the Integer GPS Ambiguities

P.J.G. Teunissen

Delft Geodetic Computing Centre (LGR)
Department of the Geodetic Engineering
Delft University of Technology
Thijssseweg 11, 2629 JA DELFT
The Netherlands

Abstract

The Global Positioning System (GPS) double-difference carrier-phase data are biased by an integer number of cycles. In this contribution a new method is introduced that enables very fast integer least-squares estimation of the ambiguities. The method makes use of an ambiguity transformation that allows one to reformulate the original ambiguity estimation problem as a new problem that is much easier to solve. The transformation aims at decorrelating the least-squares ambiguities and is based on an integer approximation of the conditional least-squares transformation. And through a flattening of the typical discontinuity in the GPS-spectrum of conditional variances of the ambiguities, the transformation returns new ambiguities that show a dramatic improvement in precision in comparison with the original double-difference ambiguities.

Contents

1	Introduction
2	Least-Squares
3	Integer Least-Squares
4	Ellipsoidal Planes of Support
5	Conditional Least-Squares
6	The Integer GPS Ambiguity Transformation
6.1	The idea of reparametrization
6.2	The admissible ambiguity transformations
6.3	On the choice of reparametrization in 2D
6.4	Bounding the triangular factor
6.5	Flattening the spectrum of conditional variances
7	Summary and concluding remarks
8	References

This contribution is based on the invited lecture given by the author, for Section IV "Theory and Methodology", at the General Meeting of the International Association of Geodesy, Beijing, China, August 1993.

1 Introduction

Nonlinear optimization, nonlinear least-squares and densities of nonlinear estimators are a trilogy of problems that are intimately related in the framework of estimation or adjustment of geodetic data. The description of physical phenomena generally proceeds through models in which a mapping, A , is defined, from a set of parameters, N , to a set of experimental outcomes, M . M is supposed to contain the image of the map A . The usual assumptions are that the map A is sufficiently smooth and that both the spaces M and N are continuous subspaces of R^m and R^n respectively. And indeed many of our geodetic estimation problems can be described as such. But not all! In particular it may happen that we know a priori that some of the parameters of interest have a discrete-like nature and are therefore not allowed to range through the whole of R . If this happens to be the case, standard gradient-like algorithms, such as for instance the Gauss-Newton method as used for solving smooth nonlinear least-squares problems fail to hold. Therefore alternative methods need to be devised for solving estimation problems in which N is of a discrete-like nature.

An important estimation problem where the parameters fail to range through the whole of R occurs in the field of GPS phase-data processing. The GPS double-difference phase ambiguities are known to be integer valued. And it is this a priori information that one would like to incorporate in the estimation procedure so as to improve the precision of the result. This is a non-trivial problem if one aims at numerical efficiency. And indeed, this topic has been a rich source of GPS-research over the last decade or so, see e.g. [2-8], [10], [17]. Starting from rather simple but time-consuming integer rounding schemes, the methods have evolved into complex and very efficient search algorithms. Nevertheless, at present times, it is still opportune to seek ways for improving the efficiency of the various search methods. This is in particular true for the real-time applications of GPS. But, to a certain extent, this is also true for some typical static applications of GPS. Because, if we are really able to significantly reduce the computational effort for estimating the integer ambiguities, it may also become worthwhile to tackle problems that have dimensions higher than the one's considered so far. And in particular this could open the way of treating all ambiguities of the various baselines in a GPS-network simultaneously.

The present study was initiated by the desire to obtain a better understanding of existing ambiguity search algorithms and through this, to obtain a better grip on the nature of the intrinsic difficulties that are associated with the problem of GPS-ambiguity fixing. And in particular, improve the computational speed of integer estimation. This would then, hopefully, also lead to ways of improving existing methods of GPS-ambiguity fixing. Based on our numerical experiences so far, we believe to have succeeded in formulating such an improved method of GPS-ambiguity fixing. Our proposal, which in part is based on existing ideas, is presented in section 5 and section 6, and is summarized in section 7.

Our proposal is based on a one-to-one transformation from the original double-difference ambiguities to a new set of integer ambiguities. And the essence of the method is that this reparametrization enables one to obtain new ambiguities which have a smaller variance and that are less correlated. The idea of transforming the double-difference ambiguities is of course not completely new. It can be recognized in the transformation from the L_1 - and L_2 - ambiguities to the wide-lane ambiguity. This transformation, however, is not one-to-one. Also this transformation is enacted at the level of a single channel and therefore does not take into account the presence of the receiver-satellite geometry. The idea of transforming the double-difference ambiguities can also be recognized to some extent in the customary practice of choosing that reference satellite which has a favourable influence on the precision of the double-difference ambiguities. In fact, as will be shown in the sequel, the one-to-one transformation from one set of double-difference ambiguities to another set having a different satellite as reference, belongs to the class of admissible transformations that will be considered in the sequel.

In order to properly judge the significance of the present contribution, it is important that one distinguishes between the following two problems of GPS-ambiguity fixing. First one has the problem of finding ways, preferably efficient ways, for fixing the ambiguities. This in fact is the *estimation* problem. Secondly, one has the problem of validating the fixed values of the ambiguities. This is the *testing* problem. And for a proper evaluation of the validity of the fixed ambiguities one will need the probability densities of the corresponding estimators; reference is made to the discussions in e.g. [2] and [13]. Although the procedures for validating the fixed ambiguities which are currently in use in practice, appear to work satisfactory, it is the author's opinion that a sound theoretical basis for these validation

procedure is still lacking. As a case in point, consider the customary practice of ambiguity validation. Usually the fixed ambiguities are validated by testing the ratio of the quadratic forms of residuals belonging to the most-likely and second most-likely integer candidates. This ratio is then tested against a critical value computed from an F-distribution. But this is incorrect, because of the stochastic dependency that exists between the two quadratic forms. Despite the importance of proper validation procedures, the present contribution only addresses the first problem and not the second. Hence, we will only be concerned with the problem of finding a numerically efficient way for estimating the integer ambiguities. There is therefore no harm in stressing, if the data is contaminated with unmodelled effects, that our method, efficient as it may be, might still come up with the wrong ambiguities.

Although the material of this contribution is intended for solving GPS-ambiguity fixing problems, we have tried, for those who are not too familiar with the typical intricacies of GPS, to refrain as much as possible from the use of standard GPS-terminology. The presentation is therefore cast in the framework of adjustment-theory and the problem of GPS-ambiguity fixing based on the least-squares criterium is formulated as an integer least-squares problem. Also detailed derivations of results are avoided in the paper. They will be taken up in a future presentation [15]. Instead we keep to the basic ideas involved and try to motivate the main results by appealing to intuition and stochastic or geometric interpretations.

This paper is organized as follows. In section 2 we briefly review the standard linear and nonlinear least-squares problem, and emphasize the (differential) geometric interpretation that can be given to least-squares problems [14].

In section 3 we introduce and define the integer least-squares problem. Again the geometric interpretation is emphasized. For reason of simplicity we assume the map $A: R^n \rightarrow R^m$ to be linear. This also allows an easier access to the more intricate details of the integer least-squares problem. It is shown how an integer least-squares problem can be decomposed into parts. This is shown both for the hybrid as well as for the non-hybrid case. The decomposition is based on the theorem of Pythagoras and allows one to solve the problem in two steps. The first step is rather straightforward and amounts to an ordinary least-squares adjustment. The second step comprises the minimization of a non-homogeneous quadratic form over

the set of integers. And it is with this second step that the intricacy of the problem manifests itself. The stepped-wise approach agrees with the approach that is usually taken in case of GPS-ambiguity fixing. Since the minimizer of the original integer least-squares problem is shown to be identical to the integer minimizer of the quadratic form, the remaining part of the sequel will focus on finding ways of solving for this integer minimizer.

The sections 4, 5 and 6 are closely related. In each one of these three sections a concept for solving integer least-squares problems is presented. The concepts of section 4 and section 5 have already been in use, in one form or another, for fixing GPS-ambiguities. The reason for including these two concepts is not only because they are of importance in their own right and that they reveal clearly the intricate nature of integer least-squares problems, but also because they pave the way for a proper discussion of the material of section 6.

The first concept is reviewed in section 4. It is based on the idea that an ellipsoidal region can be described by using the infinite set of ellipsoidal planes of support. This approach is very similar to the use of simultaneous confidence intervals in statistics for multiple comparisons [11]. Within the context of GPS-ambiguity fixing the method of Frei is based on it [6]. In this section it is shown how a finite subset of the infinite ellipsoidal planes of support can be used for selecting integer candidates from which then the sought for integer minimizer is chosen.

The second concept, which is based on the completion of squares of a quadratic form, is reviewed in section 5. This concept makes use of a triangular decomposition of the positive-definite matrix that describes the ellipsoidal region. This allows one then to come up with bounds for the integer candidates that are sharper in general than the bounds derived in the previous section. It is shown how these bounds can be employed for the formulation of an efficient search algorithm. Within the context of GPS-ambiguity fixing, examples of approaches that, in one form or another, make use of a triangular decomposition, are [2], [5] and [17]. The method of Wübbena [17], resembles the approach of the present section most. The method of Euler/Landau, [5], however, still uses an a priori chosen rectangular search window. That is, the triangular decomposition is not used for selecting integer candidates, but instead, it is used after a candidate set has been selected, for efficiently eliminating candidates.

Our proposal for efficiently solving the estimation problem of GPS-ambiguity fixing is presented in section 6. First the idea of reparametrization is introduced. It is argued that the integer least-squares problem becomes easier to solve if the reparametrization can achieve a scaling and rotation of the ellipsoidal region that will result in a region which has its principal axes (almost) parallel to the grid axes. The objective is thus, to decorrelate the least-squares ambiguities and to diagonalize their variance-covariance matrix. For GPS, this is in particular of relevance when only short timespan carrier-phase data is used. The ambiguities will then be extremely correlated, their confidence ellipsoid will be extremely elongated, and the spectrum of conditional variances of the ambiguities will then show a large discontinuity. In fact, it is the discontinuity in the spectrum of conditional variances, that prohibits an efficient search for the integer least-squares ambiguities. Although the aim is to decorrelate the ambiguities, true diagonality of the variance-covariance matrix will be difficult to reach. This is due to the fact that only a particular class of ambiguity transformations is admissible. They need to be integer and volume-preserving. Based on an integer approximation of the conditional least-squares transformation, a decorrelating two-dimensional transformation, which is both integer and volume-preserving, is introduced in subsection 6.3. It returns ambiguities with an improved precision and guarantees that the correlation coefficient stays sufficiently bounded. To tackle the n -dimensional problem, the two-dimensional transformation is repeatedly applied to pairs of conditional least-squares estimates of the ambiguities. This approach has been motivated by the presence of the typical discontinuity in the GPS spectrum of ambiguity conditional variances and is based on ideas from [9]. The success of our method is largely due to the presence of the discontinuity. And this discontinuity in the GPS-spectrum of ambiguity conditional variances also stipulates the relevance of satellite-redundancy and the use of dual-frequency data. By removing the discontinuity with the ambiguity transformation, the spectrum is flattened and lowered, and transformed ambiguities are obtained that show a dramatic improvement in precision. As a result the search for the transformed integer least-squares ambiguities can be performed in a highly efficient manner.

2 Least squares

The problem of least-squares can be formulated as the minimization problem:

$$(1) \quad \min_x \|y - A(x)\|^2, x \in R^n,$$

where: $A: R^n \rightarrow R^m$ ($m \geq n$); $\|\cdot\|^2 = (\cdot)^T Q_y^{-1}(\cdot)$ and Q_y is positive-definite. For varying values of $x \in R^n$, $A(x)$ traces (locally) an n -dimensional surface or manifold embedded in R^m . With the metric Q_y^{-1} of R^m , the scalar $\|y - A(x)\|^2$ therefore equals the distance from the datapoint y to the point $A(x)$ on the manifold. Hence, the problem of (1) corresponds to the problem of finding that point on the manifold, say $\hat{y} = A(\hat{x})$, which has least distance to y .

There are two conditions that $\hat{y} = A(\hat{x})$ needs to satisfy in order for it to be a (local) solution of (1). The first condition states that the least-squares residual vector $\hat{e} = y - \hat{y}$ should be orthogonal to the tangentspace of the manifold at the solution \hat{y} . And the second condition states that the datapoint y should lie within a hypersphere with centre \hat{y} and a radius equal to the largest principal normal curvature corresponding with the normal direction of the least-squares residual vector. Both these conditions are necessary and sufficient. And both these conditions are intrinsic in the sense that they are invariant to a change of variables or a reparametrization.

If the map $A(\cdot)$ is nonlinear (which happens to be the case in the majority of geodetic applications), the corresponding manifold traced by $A(x)$ is curved, and then generally no direct methods exist for solving (1) (there are exceptions). In this case one has to fall back on using iterative solution techniques. These iteration methods, such as the Gauss-Newton method, are usually based on repetitively solving linear or linearized least-squares problems.

If the map $A(\cdot)$ is linear, the corresponding manifold traced by $A(x)$ is flat. In this linear case the absence of curvature allows one to solve the minimization problem (1) explicitly. The corresponding linear least-squares estimates are given by the well-known formulae:

$$(2) \quad \begin{aligned} \hat{y} &= P_A y, \quad \hat{e} = P_A^\perp y, \\ \hat{x} &= A^{-1} P_A y, \quad \|\hat{e}\|^2 = \|P_A^\perp y\|^2, \end{aligned}$$

where: matrix P_A is the orthogonal projector that projects onto the range space of A and along its orthogonal complement; $P_A^\perp = I - P_A$; and A^{-1} is an (arbitrary) inverse of A . The estimates \hat{y} , \hat{e} and $\|\hat{e}\|^2$ are all unique, and the estimate \hat{x} is unique if and only if the linear map A has full rank n .

In the remaining of the sequel we will assume the map A to

be linear. This does of course restrict the focus of our discussion somewhat. But it is of importance to first understand the linear case before the nonlinear case can be tackled.

3 Integer least squares

A least-squares problem is said to be of the integer-type if the parameters are constrained to integer values. The problem of integer least-squares can therefore be formulated as the minimization problem:

$$(3) \quad \min_x \|y - Ax\|^2, x \in Z^n.$$

Compare with (1) and note that R^n has been replaced by Z^n . In order to describe the integer least-squares problem geometrically, consider the standard grid of coordinate lines in R^n . This standard grid is mapped by A into a new, but non-standard oblique grid in the data space R^m . This new grid is superimposed on the flat manifold spanned by the columnvectors of A . The set of gridpoints of this grid equals the set that follows when Z^n is mapped under A . Hence, it is from this mapped set of gridpoints in the manifold that points should be chosen as possible candidates for solving (3). In fact, the point of this mapped set of gridpoints that has the least distance to the given datapoint $y \in R^m$, minimizes the constrained objective function of (3). In the linear case there are at most 2^n such points; in the nonlinear case however, there may be even more. Once such a point has been found, the corresponding parameter vector may be obtained through an inverse mapping of A . And this solution is unique if A has full rank n . Note however that in some cases the solution may still be unique even if A is not of full rank. This happens if the set $\{x \in R^n \mid x - z - u, Au = 0\}$, where z is an integer minimizer of (3), has only one point in common with Z^n , namely z . In the following we will assume A to be of full rank.

The integer least-squares problem (3) may be decomposed in two parts. In order to show this we apply Pythagoras to get the following decomposition of the objective function of (3):

$$(4) \quad \|y - Ax\|^2 = \|P_A y - Ax\|^2 + \|P_A^\perp y\|^2.$$

The second term on the right-hand side equals the squared-norm of the least-squares residual vector \hat{e} . Since this term is independent of the parameters, the minimizer of (3) is identical to the minimizer of

$$(5) \quad \min_x \|P_A y - Ax\|^2, \quad x \in Z^n.$$

Hence, the two problems (3) and (5) are equivalent in the sense that the parameter vector that minimizes (3) also minimizes (5), and vice versa. From this follows that the solution of (3) can be obtained in two steps. In the first step an ordinary least-squares estimation is performed. This amounts to replacing y in (3) by its least-squares estimate $P_A y$. This then gives (5), which is solved in the second step. Once the minimizer of (5) has been found, the minimum value of the original objective function is obtained by adding the squared-norm of the least-squares residual vector to the minimum value of the objective function of (5).

In the above discussed integer least-squares problem the complete parameter vector x was constrained to lie in Z^n . It may also happen however that only some, but not all, of the parameters are constrained to integer values. In that case we are in a *hybrid* situation where, with $x = (x_1^*, x_2^*)^T$, we have $x_1 \in R^{n_1}$ and $x_2 \in Z^{n_2}$. The integer least-squares problem (3) is then replaced by

$$(6) \quad \min_{x_1, x_2} \|y - A_1 x_1 - A_2 x_2\|^2, \quad x_1 \in R^{n_1}, x_2 \in Z^{n_2}.$$

But also this integer least-squares problem can be decomposed into parts. First we introduce a reparametrization through the one-to-one transformation

$$(7) \quad \begin{bmatrix} \bar{x}_1 \\ x_2 \end{bmatrix} = \begin{bmatrix} I_{n_1} & (A_1^T Q_y^{-1} A_1)^{-1} A_1^T Q_y^{-1} A_2 \\ O & I_{n_2} \end{bmatrix} \begin{bmatrix} x_1 \\ x_2 \end{bmatrix}.$$

With this reparametrization we may now replace (6) by

$$(8) \quad \min_{\bar{x}_1, x_2} \|y - A_1 \bar{x}_1 - \bar{A}_2 x_2\|^2, \quad \bar{x}_1 \in R^{n_1}, x_2 \in Z^{n_2},$$

where $\bar{A}_2 = P_{A_1}^\perp A_2$. Analogous to (4), we can then decompose the objective function of (8) as

$$(9) \quad \|y - A_1 \bar{x}_1 - \bar{A}_2 x_2\|^2 = \|P_A (y - A_1 \bar{x}_1 - \bar{A}_2 x_2)\|^2 + \|P_A^\perp y\|^2.$$

Since the second term on the right-hand side is independent of the parameters we only need to consider the first term for the minimization. With $P_A = P_{A_1} P_{\bar{A}_2}$, this first term can be further decomposed as

$$(10) \quad \|P_A (y - A_1 \bar{x}_1 - \bar{A}_2 x_2)\|^2 = \|P_{A_1} (y - A_1 \bar{x}_1)\|^2 + \|P_{\bar{A}_2} (y - \bar{A}_2 x_2)\|^2.$$

Hence, it follows from (9) and (10), that (8) may be written as

$$(11) \quad \min_{\bar{x}_1 \in R^{n_1}, x_2 \in Z^{n_2}} \|y - A_1 \bar{x}_1 - \bar{A}_2 x_2\|^2 = \|P_A^\perp y\|^2 + \min_{\bar{x}_1 \in R^{n_1}} \|P_{A_1} y - A_1 \bar{x}_1\|^2 + \min_{x_2 \in Z^{n_2}} \|P_{\bar{A}_2} y - \bar{A}_2 x_2\|^2.$$

From this decomposition follows then how one can proceed to obtain the x_1 - and x_2 -minimizers of (6). First the x_2 -minimizer is obtained from solving the integer least-squares problem

$$(12) \quad \min_{x_2} \|P_{\bar{A}_2} y - \bar{A}_2 x_2\|^2, \quad x_2 \in Z^{n_2}.$$

Hence, in the hybrid case, (12) takes over the role of (5). Once the x_2 -minimizer is known, the x_1 -minimizer can be computed as follows.

Note that since the second term on the right-hand side of (11) equals zero, the corresponding minimizer is given as $\hat{\bar{x}}_1 = (A_1^T Q_y^{-1} A_1)^{-1} A_1^T Q_y^{-1} y$. This estimate together with the x_2 -minimizer allows one then to compute the x_1 -minimizer through the use of transformation (7). In the remaining of the sequel we will assume for reasons of simplicity to have a non-hybrid integer least-squares problem. Hence, we will consider (3) instead of (6).

Decomposition (4) implies that as far as the minimizer of (3) is concerned we may as well start from the minimization problem (5). If \hat{x} is the least-squares estimate of x , then $P_A y = A \hat{x}$ and (5) may be written as

$$(13) \quad \min_x \|\hat{x} - x\|_{Q_x}^2, \quad x \in Z^n,$$

where $\|\cdot\|_{Q_x}^2 = (\cdot)^T Q_x^{-1} (\cdot)$ and $Q_x^{-1} = A^T Q_y^{-1} A$.

As was pointed out already in the previous section, the minimization problem (13) may not have a unique solution. Although it is very unlikely that (13) has more than one solution, it is possible in principle that (13) has up to 2^n different minimizers. Still however, we will assume in the present sequel that (13) has one and only one solution. Our motivation for this assumption stems from the way the integer least-squares problem is applied in the context of GPS-ambiguity fixing. If x stands for the vector of double-difference ambiguities, a non-unique solution for (13) will namely imply that a reliable fixing of the ambiguities is not feasible. In the remaining of the sequel we will focus on solving (13).

Unfortunately there are in general no standard techniques available for solving (13) as they are available for solving ordinary least-squares problems. As a consequence one has to resort to methods that in one way or another make use of

a discrete search strategy for finding the minimizer of (13). The idea is to use the objective function of (13) for introducing an ellipsoidal region in R^n , on the basis of which a search is performed for the minimizer of (13). The ellipsoidal region is given by

$$(14) \quad (x-\hat{x})^T Q_{\hat{x}}^{-1} (x-\hat{x}) \leq \chi^2.$$

Through the selection of the positive constant χ^2 the size of the ellipsoidal region can be controlled. However, already with the selection of χ^2 care has to be excersized. A small value for χ^2 may result in an ellipsoidal region that fails to contain the minimizer of (13). And a too large value for χ^2 may result in a region for which the search for the minimizer becomes too time-consuming. Unfortunately it is difficult to formulate a data-independent criterion for selecting χ^2 , that ensures that the region contains one or more gridpoints. This is due to the fact that the ellipsoidal region is centered at $\hat{x} \in R^n$ and not centered at a gridpoint of Z^n . Would the latter be the case, then the volume of the ellipsoid could be used to set a reference value for χ^2 . It can namely be shown that for the case $\hat{x} \in Z^n$, the ellipsoidal region (14) would at least contain one gridpoint other than \hat{x} if its volume is larger than or equal to 2^n (which is the volume of the cube $|x_i| \leq 1, i=1, \dots, n$). In our case an alternative approach has to be taken. One approach would be to round all the individual coordinates of \hat{x} to their nearest integer, substitute the so obtained vector for x in the left-hand side of (14) and then take χ^2 to be equal to the function value of the quadratic form. This approach at least ensures that (14) contains minimally one gridpoint. However, the so obtained value for χ^2 may also be overly conservative. And this may occur especially when the ellipsoidal region is extremely elongated (which is typically the case with GPS when the observational timespan is short). In the context of selecting χ^2 , it is of interest to note that our numerical experiments indicate that the *volume* of the ellipsoidal region (14) gives a good approximation to the number of integer vectors that lie within the ellipsoidal region. This suggests that one could use the volume of the ellipsoidal region as indicator to decide whether or not the scalar χ^2 should be scaled down or scaled up.

An alternative and from a statistical testing point of view more appealing approach would be to rely on and to make use of the statistical distribution of the least-squares estimator of x . If the observables are normally distributed with mean Ax and variance matrix Q_y , then the quadratic form of (14) has a central Chi-square distribution with $m-n$ degrees of

freedom (it is a central F -distribution if $Q_{\hat{x}}$ has been scaled with the a posteriori variance factor). As reference value for χ^2 one may now choose χ^2 to be equal to the α -percentage point of the Chi-square (or F -) distribution. With this choice for χ^2 one is of course not certain that (14) indeed contains a gridpoint. But the choice does ensure that (14) contains a grid point with probability $1-\alpha$. This gridpoint is the mean of the least-squares estimator of x .

The above shows that one should give some consideration to the way χ^2 is selected. If the only objective is to solve the minimization problem (13), then χ^2 should be chosen in a way that guarantees that (14) contains the minimizer. If however the objective is also to statistically validate the minimizer, then the approach based on the α -percentage point of the Chi-square (or F -) distribution can be used. Because, if α is chosen small enough and (14) still fails to contain a grid point, then the minimizer of (13) can be considered to be invalidated.

From now on it will be assumed that a value for χ^2 has been selected. With this value for χ^2 the ellipsoidal region (14) is then taken as the point of departure for developing a search strategy to obtain the minimizer of (13). Different search strategies are possible and some of them have in fact been proposed already in the GPS-literature. In the following we will review two of such concepts that already have been in use for GPS ambiguity fixing. They are based on using the planes of support of an ellipsoid and on completing a quadratic equation to squares. These two concepts are reviewed in sections 4 and 5.

4 Ellipsoidal planes of support

One way of finding the minimizer of (13) is to identify first the set of gridpoints that satisfy the inequality (14) and then to pick that gridpoint that gives the smallest function value for the quadratic form. However, the quadratic form of (14) can not be used as such to identify the set of candidate gridpoints. The idea is therefore to replace inequality (14) with an equivalent description that is based on using the planes of support of the ellipsoid. This equivalence can be constructed as follows: Let a be an arbitrary vector of R^n and let $x-\hat{x}$ be orthogonally projected onto a . The orthogonal projection (where orthogonality is measured with respect to the metric $Q_{\hat{x}}$) of $x-\hat{x}$ onto a is then given as: $a(a^* Q_{\hat{x}}^{-1} a)^{-1} a^* Q_{\hat{x}}^{-1} (x-\hat{x})$. And the square of the length of this vector reads: $[a^* Q_{\hat{x}}^{-1} (x-\hat{x})]^2 / (a^* Q_{\hat{x}}^{-1} a)$. Now, since the length of the orthogonal projection of a vector onto an arbitrary direction is always less than or equal to the length of the vector

itself, we have

$$(x-\hat{x})^T Q_{\hat{x}}^{-1} (x-\hat{x}) = \max_{a \in R^n} \frac{[a^T Q_{\hat{x}}^{-1} (x-\hat{x})]^2}{a^T Q_{\hat{x}}^{-1} a}$$

From this follows then, when a is replaced by $Q_{\hat{x}} c$, the equivalence

$$(15) \quad (x-\hat{x})^T Q_{\hat{x}}^{-1} (x-\hat{x}) \leq \chi^2 \Leftrightarrow \frac{[c^T (x-\hat{x})]^2}{c^T Q_{\hat{x}} c} \leq \chi^2, \forall c \in R^n.$$

Both type of inequalities describe the same ellipsoidal region. In the second type we recognize $c^* (x-\hat{x}) = \pm (c^* Q_{\hat{x}} c)^{1/2} \chi$, which is the pair of parallel planes of support of the ellipsoid having vector c as normal. The above equivalence therefore states that the ellipsoidal region coincides with the region that follows from taking all intersections of the areas between each pair of ellipsoidal planes of support. Hence, in order to find the candidate gridpoints that satisfy (14) we may as well make use of the ellipsoidal planes of support.

When working with the above equivalence for our purposes, there are however two restrictions that need to be appreciated. First of all, the above equivalence only holds for the *infinite* set of planes of support. But for all practical purposes one can only work with a finite set. Working with a finite set implies however that the region bounded by the planes of support will be larger in size than the original ellipsoidal region. Of course, one could think of minimizing the increase in size by choosing an appropriate set of normal vectors c . For instance, if the normal vectors c are chosen to be in the direction of the major and minor axes of the ellipsoidal region, then the resulting region will fit the ellipsoid best. But here is where the second restriction comes into play. One simply has no complete freedom in choosing the planes of support. Their normals c should namely be chosen such that the resulting interval $[c^* (x-\hat{x})]^2 \leq c^* Q_{\hat{x}} c \chi^2$ can indeed be used for selecting candidate grid points. Hence, the normals c can not be chosen arbitrarily.

The simplest approach to the above would be to submit oneself to this situation and to choose the normals to be parallel to the grid axes. When the normals are chosen as $c_i = (0, \dots, 1, 0, \dots, 0)^*$, with the 1 as the i th-coordinate, the region bounded by the planes of support becomes

$$(16) \quad (x_i - \hat{x}_i)^2 \leq \sigma_{\hat{x}_i}^2 \chi^2, \quad i=1, \dots, n,$$

where $\sigma_{\hat{x}_i}^2$ is the variance of the least-squares estimator of x_i . The intervals of (16) can be used to select candidate gridpoints from which then the minimizer of (13) can be chosen.

Although the approach based on (16) is certainly a valuable one, it will be clear that it can become quite time-consuming when the region defined by (16) is significantly larger than the original ellipsoidal region. And this will definitely be the case when the ellipsoid is both elongated and rotated with respect to the grid axes. Of course, one can reduce the size of the region by introducing additional planes of support. For instance, in addition to $c_i = (0, \dots, 1, 0, \dots, 0)^*$, admissible choices for the normals are the sum or differences of c_i and c_j , $i, j=1, \dots, n$, $i \neq j$. And depending on the elongation and orientation of the ellipsoid, this may indeed significantly reduce the size of the region enclosed by the planes of support. The problem remains however that this way of including additional planes of support is somewhat ad hoc and need not necessarily lead to a significant reduction in size of the region.

One conclusion that can be drawn from the above discussion is that the efficacy of the method depends to a large extent on the elongation and orientation of the ellipsoid. This observation has been the motivation for developing the method that will be discussed in section 6. First however we will review another interesting concept that already has been in use, in one form or another, for GPS ambiguity fixing. This concept is based on completing the square of a quadratic equation, or phrased in statistical terms, it is based on a sequential conditional least-squares adjustment.

5 Conditional Least-Squares

When use is made of the planes of support as previously discussed, all bounds (such as $\sigma_{\hat{x}_i}^2 \chi^2$ in (16)) are set prior to the actual search. That is, the setting of the bounds is independent of the search process. One may wonder however whether it is not possible to keep adjusting these bounds during the search process. For instance, let x be partitioned as $x = (x_1^*, x_2^*)^*$ and assume that already candidate integers have been found for the elements of x_1 . Then clearly it is possible to formulate new bounds for the elements of x_2 that are sharper in general than the original bounds in the previous section. And as it turns out, this idea can be implemented very efficiently when use is made of completing the square of a quadratic equation.

In order to describe the method we will start with the two-dimensional case first. Let the two-dimensional ellipsoidal region be given as

$$(17) \quad ax_1^2 + 2bx_1x_2 - cx_2^2 \leq \chi^2,$$

where $a>0, c>0, ac-b^2>0$. For the moment we simply assume the ellipse to be centered at the origin. When we use the approach of the previous section and apply (16), the two bounds for x_1 and x_2 follow as

$$(18) \quad x_1^2 \leq \begin{bmatrix} 1 \\ 0 \end{bmatrix} \begin{bmatrix} a & b \\ b & c \end{bmatrix}^{-1} \begin{bmatrix} 1 \\ 0 \end{bmatrix} \chi^2 = \chi^2 / (a - b^2/c),$$

and

$$(19) \quad x_2^2 \leq \begin{bmatrix} 0 \\ 1 \end{bmatrix} \begin{bmatrix} a & b \\ b & c \end{bmatrix}^{-1} \begin{bmatrix} 0 \\ 1 \end{bmatrix} \chi^2 = \chi^2 / (c - b^2/a).$$

Bound (18) is the sharpest possible bound for x_1 when nothing is known about x_2 . Similarly, bound (19) is the sharpest possible bound for x_2 when nothing is known about x_1 . But each of these bounds can be improved once the other parameter is known. This can be seen when (17) is completed to a sum of squares. Completing the square of (17) gives

$$(20) \quad a(x_1 - \frac{b}{a}x_2)^2 - (c - \frac{b^2}{a})x_2^2 \leq \chi^2.$$

And from this follows that one can bound x_1 as

$$(21) \quad (x_1 - \frac{b}{a}x_2)^2 \leq \frac{c - b^2/a}{a} \left[\frac{\chi^2}{c - b^2/a} - x_2^2 \right].$$

This shows, since $a \geq a - b^2/c$ and $x_2^2 \geq 0$, that except for the trivial case $b=0$, the bound of (21) is always sharper than the bound of (18). Advantage can therefore be gained from replacing (18) by (21). And this gain may be considerable when the ellipse is elongated and rotated with respect to the grid axes. Since the bound of (21) depends on x_2 , first (19) should be used to come up with an integer candidate for x_2 . This step is then followed by (21) for determining an integer candidate for x_1 . Once a pair of integer candidates has been found, a new and smaller reference value for χ^2 can be computed. This enables us then to shrink the ellipsoidal region and to perform a renewed search for integer candidates. In this way one can efficiently scan the ellipsoidal region to locate the sought for minimizer of (13).

In order to generalize the above to the multi-dimensional case, we first note that (20) can be written in vector-matrix form as

$$(22) \quad \begin{bmatrix} x_1 \\ x_2 \end{bmatrix} \begin{bmatrix} 1 & 0 \\ b/a & 1 \end{bmatrix} \begin{bmatrix} a & 0 \\ 0 & c - b^2/a \end{bmatrix} \begin{bmatrix} 1 & 0 \\ b/a & 1 \end{bmatrix} \begin{bmatrix} x_1 \\ x_2 \end{bmatrix} \leq \chi^2.$$

This shows that completing the square corresponds to a triangular decomposition. Also note that with

$$\begin{bmatrix} a & b \\ b & c \end{bmatrix} = \begin{bmatrix} \sigma_1^2 & \sigma_{12} \\ \sigma_{21} & \sigma_2^2 \end{bmatrix}^{-1},$$

$$\begin{bmatrix} 1 & 0 \\ b/a & 1 \end{bmatrix} = \begin{bmatrix} 1 & -\sigma_{12}\sigma_2^{-2} \\ 0 & 1 \end{bmatrix}, \text{ and}$$

$$\begin{bmatrix} a & 0 \\ 0 & c - b^2/a \end{bmatrix} = \begin{bmatrix} (\sigma_1^2 - \sigma_{12}^2/\sigma_2^2)^{-1} & 0 \\ 0 & 1/\sigma_2^2 \end{bmatrix},$$

the inequality of (20) can be written as

$$(23) \quad (\hat{x}_{1|2}/\sigma_{1|2})^2 + (x_2/\sigma_2)^2 \leq \chi^2,$$

and the two inequalities (19) and (21) can be written as

$$(24) \quad x_2^2 \leq \sigma_2^2 \chi^2 \text{ and } \hat{x}_{1|2}^2 \leq \frac{\sigma_{1|2}^2}{\sigma_2^2} (\sigma_2^2 \chi^2 - x_2^2),$$

in which we recognize the conditional least-squares estimate $\hat{x}_{1|2} = x_1 - \sigma_{12}\sigma_2^{-2}x_2$ and its variance $\sigma_{1|2}^2 = \sigma_1^2 - \sigma_{12}^2/\sigma_2^2$. This shows that the triangular decomposition corresponds to a conditional least-squares adjustment [1], [12]. To generalize we therefore will use a *sequential conditional least-squares adjustment* for the multi-dimensional case.

A sequential conditional least-squares adjustment (or an *LDU*-decomposition of $Q_{\hat{x}}^{-1}$), starting with a conditioning on x_n and ending with a conditioning on x_1 , will then in analogy to (23) give for the multi-dimensional case,

$$(25) \quad \sum_{i=1}^n (x_i - \hat{x}_{i|i-1, \dots, n})^2 / \sigma_{i|i-1, \dots, n}^2 \leq \chi^2.$$

From this and in analogy to (24), we can construct the following bounds

$$(26) \quad (x_i - \hat{x}_{i|i-1, \dots, n})^2 \leq \lambda_{\hat{x}_i} \sigma_{i|i-1, \dots, n}^2 \chi^2, \quad i = 1, \dots, n,$$

with

$$(27) \quad \lambda_{\hat{x}_i} = \left[1 - \sum_{j=i+1}^n \frac{(x_j - \hat{x}_{jj|i-1, \dots, n})^2}{\sigma_{jj|i-1, \dots, n}^2 \chi^2} \right]$$

Here we have made use of the notation $\hat{x}_{i|i-1, \dots, n}$ to denote the least-squares estimate of x_i conditioned on fixing $x_j, j=i+1, \dots, n$.

Compare (26) with (16). Since clearly $0 \leq \lambda_i \leq 1$, and since a conditional variance is always smaller than or equal to its

unconditional counterpart, $\sigma_{i|i-1, \dots, n}^2 \leq \sigma_{x_i}^2$, it follows that the bounds of (26) (or (23)) are always sharper than or at least as sharp as the bounds of (16).

Also note the regularity in the above bounds. The bound for x_i is equal to the gap in the previous bound times the ratio $\sigma_{i|i-1, \dots, n}^2 / \sigma_{i-1|i-2, \dots, n}^2$. The above set of inequalities can now be used as follows for solving (13). First a candidate integer is determined for x_n using the first bound. This candidate integer is then used to compute the second bound, from which a candidate integer for x_{n-1} is determined. This process is continued up to the point that a complete vector x with candidate integer coordinates is constructed. With this vector one can then shrink the ellipsoidal search region and perform a renewed search within the shrunken ellipsoidal region. Repeated application of the above steps will then finally lead to the minimizer of (13). Note that the renewed search need not necessarily commence with x_n . If for instance it is known that due to the shrinking of χ^2 , the integers $x_n, x_{n-1}, \dots, x_{k+1}$ are the only candidates, the renewed search may commence with x_k .

It may happen that the above procedure halts before a complete candidate vector x has been found. This occurs when the size of the bound is such that no candidate integer lies within the interval. If this happens to be the case, one should return to the previous bound and increase (or decrease) the previously found candidate integer by one. From there on one can then continue again.

It can be deduced that the bounds of (26) have the *tendency* to become smaller as the index i gets smaller. Clearly λ_i gets smaller as the index i gets smaller. But generally also $\sigma_{i|i-1, \dots, n}^2$ has the tendency to get smaller as the index i gets smaller. The more constraints are included the smaller the conditional variance gets. This tendency can also be explained if we interpret $\sigma_{i|i-1, \dots, n}^2$ geometrically. It can be shown that

$$\sigma_{i|i-1, \dots, n}^2 = |P_{A_{(i)}}^\perp a_i|^2 = |a_i|^2 \sin^2 \alpha_i$$

where: a_i is the i th-column vector of matrix A , $A_{(i)}$ is the matrix that follows from taking the first $(i-1)$ -number of column vectors from A and α_i is the angle between vector a_i and the range space of matrix $A_{(i)}$. Now, the angle α_i will have the tendency to become smaller the larger the dimension of the range space of $A_{(i)}$ gets. In fact the angle will be zero when the dimension of the range space equals n . Therefore, when the lengths of the column vectors of matrix A are approximately constant, also $\sigma_{i|i-1, \dots, n}^2$ will have the tendency

to decrease when the index i gets larger.

The above suggests that in order to reduce the potential of halting, it may be worthwhile to order the elements of \hat{x} according to their (conditional) precision. Because, even if the bound of (26) is small at a certain level $i=l$, halting will not take place as long as $\hat{x}_{l|l-1, \dots, n}$ is sufficiently pushed towards an integer value. But this requires that the previously chosen integers $x_i, i=l+1, \dots, n$, are indeed coordinates of an integer vector x that lies within the ellipsoidal region. And the probability that this will be the case is higher the better the precision of these elements is.

The problem that the search for the integer candidate vector halts, is a serious one in case of GPS carrier-phase processing. For a single baseline model, it can namely be shown, see [15], that the spectrum of conditional variances of the ambiguities, $\sigma_{i|i-1, \dots, n}^2$ for $i=n, \dots, 1$, has a large discontinuity when passing from $\sigma_{n-2|n-1, n}^2$ to $\sigma_{n-3|n-2, n-1, n}^2$. In fact, one can show that $\sigma_{j|j-1, \dots, n}^2 \ll \sigma_{i|i-1, \dots, n}^2$ for $j=1, \dots, n-3$ and $i=n-2, n-1, n$. This implies, since the first three bounds of (26) will be rather loose, that quite a number of integer-triples satisfy these bounds. But, this on its turn implies, when we start working with the fourth bound, which is very tight due to the steep decrease in value of the conditional variances, that we have a high likelihood of not being able to find an integer candidate that satisfies this fourth bound. The potential of halting is therefore very significant when one passes from the third to the fourth bound. As a consequence a large number of trials are required, before one is able to move on to the next bound. And it is this inefficiency, that will be tackled by our method proposed in the next section. The method that will be introduced in the next section, overcomes the problem of halting, through a flattening and a lowering of the level of the GPS spectrum of ambiguity conditional variances.

6 The integer GPS ambiguity transformation

6.1 The idea of reparametrization

In the previous two sections we have dealt with two ways of solving the integer least-squares problem (5) (or (13)). First the use of the ellipsoidal planes of support was discussed. But as was pointed out, the bounds that follow from using the ellipsoidal planes of support can be rather conservative, in particular when the ellipsoid is elongated and rotated with respect to the grid axes. Moreover, these bounds are fixed from the outset. This observation then led

to the idea to introduce adjustable bounds, bounds that are made dependent on the stage of progress of the search process. These bounds were obtained through a sequential conditional least-squares adjustment, which resulted in the introduction of the conditional least-squares estimates $\hat{x}_{i|1,\dots,n}$, $i=1,\dots,n$. And it was shown that these bounds are indeed much less conservative. Up to this point however, we have been working solely on the basis of representation (14). But one may wonder whether it is not possible to obtain a further improvement in the search process, if one can replace (14) with an alternative but equivalent representation. And this indeed turns out to be the case. A new idea of the present section is therefore to *reparametrize* the integer least-squares problem such that an equivalent formulation is obtained, but one that is much easier and hence much faster to solve. In order to understand what our reparametrization should achieve, we first pause for a moment to present two ways of visualizing the integer least-squares problem. We will start with the data space point of view.

Assume that $Q_y = I_m$ (if this is not the case one simply has to replace in the following, A by $Q_y^{-1/2}A$ and $P_A y$ by $Q_y^{-1/2}P_A y$), and let $G_y = \{y \in R^m | y = Ax, x \in Z^n\}$ be the set of gridpoints that is generated by the column vectors of A . Then (5) amounts to finding that element of G_y that has the least (cartesian) distance to the least-squares estimate $P_A y$. In general this is a nontrivial problem to solve. The intricacy of the problem stems from the fact that although the metric is standard ($Q_y = I_m$), the grid G_y is not (the columnvectors of A are oblique in general). The problem becomes trivial however if in addition to the metric being standard, also the grid is standard (or at least orthogonal). Because if this happens to be the case, then the columnvectors of A are mutual orthogonal and (5) can simply be solved as follows.

First the consistent system of equations

$$P_A y = Ax$$

is solved for x , giving the least-squares estimate \hat{x} . And then the minimizer of (5) is obtained from a simple rounding of the individual elements of \hat{x} to the nearest integer.

The same conclusion is reached if we visualize the integer least-squares problem from a parameter space point of view. But contrary to the data space point of view we now have a non-standard metric with a standard grid. In formulation (13), x ranges namely through the standard set of gridpoints of R^n , which is Z^n , whereas distance is now measured with a non-standard metric, namely Q_x^{-1} . But as with the data space point of view we again observe that (13) becomes a trivial problem once both the metric and grid are standard. The

conclusion reads therefore that the integer least-squares problem (5) or ((13)) can simply be solved by means of rounding, if the columnvectors of A are mutual orthogonal, or if the matrix Q_x is diagonal. In order to see what happens to the methods of section 4 and 5, when Q_x is diagonal, consider the following. If Q_x is diagonal, the orientation of the ellipsoid is such that its major and minor axes are parallel to the grid axes. And in that case the n -dimensional rectangular box defined by (16) indeed fits the ellipsoid best. In that case the conditional variances also reduce to ordinary variances, and the left-hand sides of (23) become identical to those of (16). The bounds of (23) remain however sharper than those of (16).

Since the situation where Q_x is diagonal is the best one can hope for in any integer least-squares problem, we will try to find ways to come as close as possible to this ideal situation. This in short, is the essence of the method of this section.

6.2. The admissible ambiguity transformations

Let Z be an n -by- n matrix of full rank and define

$$(28) \quad z = Z \cdot x, \quad \hat{z} = Z \cdot \hat{x}, \quad Q_z = Z \cdot Q_x Z.$$

Then

$$(29) \quad (x - \hat{x}) \cdot Q_x^{-1} (x - \hat{x}) = (z - \hat{z}) \cdot Q_z^{-1} (z - \hat{z}).$$

The variance-covariance matrix Q_z is clearly diagonal if matrix Z contains the eigenvectors of Q_x . Unfortunately however, this choice for matrix Z is not admissible in case of our *integer* least-squares problem. Because, if this choice for Z would be used, the vector in Z^n that follows from rounding the coordinates of \hat{z} to their nearest integer would in general fail to produce a vector x in Z^n . In other words, if $z \in Z^n$ then generally $x = Z^{-1} \cdot z \notin Z^n$. This dilemma points out that only a restricted class of transformations qualifies for reparametrizing the integer least-squares problem. Fortunately, this class of transformations can easily be characterized [16]. They need to be *volume preserving* and have elements which are *integers*. Typical examples of matrices that fall in this class are the identity-matrix and the permutation matrices. Note that we already made use of the permutation matrices in the previous section when ordering the elements of \hat{x} according to their (conditional) precision. Also within the context of GPS double-difference ambiguity fixing, one in fact already has been using transformations that are volume preserving and have integer elements. This is the case when one changes from one set of double-

difference ambiguities to another set having a different satellite as reference. This can be seen as follows. Take as an example the situation that five satellites are available. The single-difference ambiguity related to satellite i is denoted as a_i , and the corresponding double-difference ambiguity having satellite j as reference is denoted as $a_i^{(j)} = a_i - a_j$. The regular transformation from $a_i^{(1)}$ to $a_i^{(2)}$ reads then

$$\begin{pmatrix} a_1^{(2)} \\ a_3^{(2)} \\ a_4^{(2)} \\ a_5^{(2)} \end{pmatrix} = \begin{pmatrix} -1 & 0 & 0 & 0 \\ -1 & 1 & 0 & 0 \\ -1 & 0 & 1 & 0 \\ -1 & 0 & 0 & 1 \end{pmatrix} \begin{pmatrix} a_2^{(1)} \\ a_3^{(1)} \\ a_4^{(1)} \\ a_5^{(1)} \end{pmatrix}.$$

The transformation matrix in this expression has integer elements and its determinant equals -1. Hence, it follows that this matrix is indeed a member of the class of admissible ambiguity transformations. In a similar way one can show that all transformations that change the reference satellite of the double-difference ambiguities belong to the above mentioned class of admissible transformations.

By choosing a matrix Z from the above mentioned class, we can now replace our original integer least-squares problem (13) by the equivalent, but *reparametrized* integer least-squares problem

$$(30) \quad \min_z (z - \hat{z})^T Q_z^{-1} (z - \hat{z}), \quad z \in Z^n.$$

And once the minimizer of (30) has been found, the minimizer of (13) can be recovered from invoking $x = Z^{-1}z$.

It will be clear that because of the restrictions on Z , no true diagonality of Q_z can be hoped for. This leaves us with two questions. Firstly, how to measure the non-diagonality of Q_z , and secondly, how to choose matrix Z so as to obtain near-diagonality. In order to answer the first question, we note that the variance-covariance matrix Q_z is diagonal if and only if its *correlation* matrix R_z equals the identity matrix. That is, Q_z is diagonal if and only if all elements of \hat{z} are fully decorrelated. In the two-dimensional case the determinant of R_z is related to the correlation coefficient as: $\det(R_z) = 1 - \rho_z^2$. This shows that for the two-dimensional case, Q_z is diagonal if and only if $\det(R_z) = 1$. But, it can be shown that this also holds true for dimensions higher than two [15]. We therefore introduce as measure of diagonality of Q_z , the scalar

$$(31) \quad r_z = \det(R_z)^{\frac{1}{2}} \quad (0 \leq r_z \leq 1).$$

Since the scalar r_z measures the decorrelation between the

elements of \hat{z} , it will be referred to as the *decorrelation number* of \hat{z} . The elements of \hat{z} are fully decorrelated when r_z equals one, and they are poorly decorrelated when r_z is close to zero. It follows from the triangular decomposition of Q_z , that r_z is related to the spectrum of *conditional* and *unconditional* standard deviations as

$$(32) \quad r_z = \prod_{i=1}^n \frac{\sigma_{i|i-1, \dots, n}}{\sigma_{z_i}}.$$

From the fact that the nominator in this expression is independent of Z , since $\det(Q_z) = \det(Z^T Q_x Z) = \det(Q_x)$, follows, that the elements of \hat{z} are less correlated than those of \hat{x} , $r_z > r_x$, when $\sigma_{x_i} > \sigma_{z_i}$. Hence, the variance-covariance matrix Q_z is less *non-diagonal* than Q_x , when its *diagonal* elements are smaller than those of Q_x . The gain in decorrelation can be measured by the ratio

$$(33) \quad r_z / r_x = \prod_{i=1}^n \sigma_{x_i} / \sigma_{z_i}.$$

This gain can be given the following geometrical interpretation. The volume of the n -dimensional rectangular box (16) that encloses the ellipsoidal region (14) is given as $2^n \chi^n \prod \sigma_{x_i}$. Similarly, the n -dimensional box that encloses the ellipsoidal region defined by Q_z , has volume $2^n \chi^n \prod \sigma_{z_i}$. This shows, that it is the relative decrease in volume of the n -dimensional box, that directly measures the gain in decorrelation. A maximum decrease in volume is achieved when $\prod_{i=1}^n \sigma_{z_i} = \prod_{i=1}^n \sigma_{x_i}$, in which case $r_z = 1$.

6.3 On the choice of reparametrization in 2D

In order to answer the question as to how to construct matrix Z , we first consider the problem in two dimensions. Let \hat{x} and Q_x be given as

$$(34) \quad \hat{x} = \begin{pmatrix} \hat{x}_1 \\ \hat{x}_2 \end{pmatrix} \text{ and } Q_x = \begin{pmatrix} \sigma_1^2 & \sigma_{12} \\ \sigma_{21} & \sigma_2^2 \end{pmatrix}.$$

We now need to come up with a matrix Z , which has integer entries, which is volume-preserving, and which allows us to decorrelate the two elements of \hat{x} . We already know from section 5, see equation (23), that a complete decorrelation is obtained, when \hat{x}_1 is replaced by its corresponding *conditional least-squares* estimate $\hat{x}_{1|2}$. The transformation that achieves this, is given as

$$(35) \quad \begin{pmatrix} \hat{x}_{1|2} \\ \hat{x}_2 \end{pmatrix} = \begin{pmatrix} 1 & -\sigma_{12}\sigma_2^{-2} \\ 0 & 1 \end{pmatrix} \begin{pmatrix} \hat{x}_1 \\ \hat{x}_2 \end{pmatrix}$$

For reasons of convenience, we have assumed the expectations of \hat{x}_1 and \hat{x}_2 to be simply zero for the moment. Note, that transformation (35) not only decorrelates, but, in line with the correspondence between linear least-squares estimation and best linear unbiased estimation, also returns $\hat{x}_{1|2}$ as the element which has the best precision of all linear unbiased functions of \hat{x}_1 and \hat{x}_2 . Instead of using (35), we can ofcourse also interchange the role of the two entries of \hat{x} and use the transformation

$$(36) \quad \begin{pmatrix} \hat{x}_1 \\ \hat{x}_{2|1} \end{pmatrix} = \begin{pmatrix} 1 & 0 \\ -\sigma_{21}\sigma_1^{-2} & 1 \end{pmatrix} \begin{pmatrix} \hat{x}_1 \\ \hat{x}_2 \end{pmatrix}.$$

Both of the above transformations fully decorrelate. Also note, that both transformations are volume-preserving. Hence, the only condition that is still not satisfied, is the condition that the entries of the transformation need to be integer. In order to repair this situation, we *approximate* the above transformations by replacing $\sigma_{21}\sigma_1^{-2}$ by $[\sigma_{21}\sigma_1^{-2}]$, or $\sigma_{12}\sigma_2^{-2}$ by $[\sigma_{12}\sigma_2^{-2}]$, where $[.]$ stands for "rounding to the nearest integer". The volume-preserving property is retained by this integer-approximation. The decorrelation-property is however, not retained. But still, one can show, although the integer-approximation does not allow for a complete decorrelation, that it allows one to improve the precision of the elements and that it allows one to bound the correlation between the elements when the two transformations are used in an alternating fashion. Based on (35) and (36), the idea is therefore to use the following two type of transformations

$$(37) \quad Z_1^* = \begin{pmatrix} 1 & z_{12} \\ 0 & 1 \end{pmatrix} \text{ and } Z_2^* = \begin{pmatrix} 1 & 0 \\ z_{21} & 1 \end{pmatrix}$$

in which z_{12} and z_{21} are appropriately chosen integers. The two type of transformations are applied in such a way, that they replace the element with the poorest precision with one that has an improved precision. Thus, when $\sigma_2^2 \leq \sigma_1^2$, we start with Z_1^* and we choose the scalar z_{12} as $z_{12} = -[\sigma_{12}\sigma_2^{-2}]$. This gives

$$(38) \quad Z_1^* Q_x Z_1^* = \begin{pmatrix} \sigma_1^2 & \sigma_{1'2} \\ \sigma_{2'1} & \sigma_2^2 \end{pmatrix} \text{ with } \sigma_1^2 \leq \sigma_1^2.$$

Then, if the precision of the first element is still not better than that of the second, $\sigma_2^2 \leq \sigma_1^2$, we stop, else we continue with Z_2^* and choose z_{21} as $z_{21} = -[\sigma_{21}\sigma_1^{-2}]$. This gives then

$$(39) \quad Z_2^* Z_1^* Q_x Z_1 Z_2 = \begin{pmatrix} \sigma_1^2 & \sigma_{1'2} \\ \sigma_{2'1} & \sigma_2^2 \end{pmatrix} \text{ with } \sigma_2^2 \leq \sigma_2^2.$$

Then, if the precision of the second element is still not better than that of the first, $\sigma_1^2 \leq \sigma_2^2$, we stop, else we continue again with Z_1^* and choose z_{12} as $z_{12} = -[\sigma_{1'2}\sigma_2^{-2}]$. This whole process of alternatingly using Z_1^* and Z_2^* finally stops when one fails to improve the precision of the elements. And when this happens, the correlation coefficient is bounded as $\rho_{\hat{x}}^2 \leq \frac{1}{4}$, since then both of the inequalities, $|\sigma_{1'2}\sigma_2^{-2}| \leq \frac{1}{2}$ and $|\sigma_{2'1}\sigma_1^{-2}| \leq \frac{1}{2}$, are satisfied.

Geometrically, the above sequence of transformations can be given the following useful interpretation. Imagine the confidence-ellipse of \hat{x} . The first transformation Z_1^* then pushes the two *vertical* tangents of the ellipse from the $\pm\sigma_1\chi$ level towards the $\pm\sigma_1\chi$ level, while at the same time keeping fixed the volume (area) of the ellipse and the location of the two horizontal tangents of the ellipse. The second transformation Z_2^* then pushes the two *horizontal* tangents of the ellipse from the $\pm\sigma_2\chi$ level towards the $\pm\sigma_2\chi$ level, while at the same time keeping fixed the volume of the ellipse and the location of the two vertical tangents. And this process is continued until the next transformation reduces to the trivial identity. Since the volume of the ellipse is kept constant at all times, whereas the volume of the enclosing rectangular box is reduced in each step, it follows that not only the decorrelation number gets improved, but also that the shape of the ellipse is forced to become more sphere-like.

Once the above sequence of transformations that make up Z^* has been applied, we have $\rho_{\hat{x}}^2 \leq \frac{1}{4}$, which implies for the decorrelation number that

$$(40) \quad r_{\hat{x}}^2 \geq 3/4.$$

This is a very significant result, since we know that the original double-difference ambiguities are extremely correlated when based on short timespan carrier-phase data. From this bound also follows, together with $\sigma_{\hat{x}_2}^2 \leq \sigma_{\hat{x}_1}^2$ and $\sigma_{\hat{x}_{1|2}}^2 = r_{\hat{x}}^2 \sigma_{\hat{x}_1}^2$, that

$$(41) \quad \sigma_{\hat{x}_{1|2}}^2 \geq \frac{3}{4} \sigma_{\hat{x}_2}^2.$$

Hence, the transformation Z^* guarantees that the transformed conditional variance $\sigma_{\hat{x}_{1|2}}^2$ will never be much smaller than

the variance $\sigma_{\hat{x}_2}^2$. But this implies, that the transformation removes to a large extent any discontinuity that might be present in the original variances, $\sigma_{\hat{x}_1}^2 \ll \sigma_{\hat{x}_2}^2$. And as we observed earlier in section 5, this is precisely the situation that we are confronted with in case of GPS carrier-phase processing.

6.4 Bounding the triangular factor

In order to obtain a higher-dimensional version of Z^* , we first try to find a generalization of Z_1^* . In two dimensions, \hat{x} is transformed by Z_1^* as

$$(42) \quad \begin{pmatrix} \hat{x}_1 \\ \hat{x}_2 \end{pmatrix} = \begin{pmatrix} 1 & -[\sigma_{12}\sigma_2^{-2}] \\ 0 & 1 \end{pmatrix} \begin{pmatrix} \hat{x}_1 \\ \hat{x}_2 \end{pmatrix}$$

If we substitute the inverse of (35) into the right-hand side of (42) and apply the error propagation law, we get

$$(43) \quad \begin{pmatrix} \sigma_{\hat{x}_1}^2 & \sigma_{\hat{x}_1\hat{x}_2} \\ \sigma_{\hat{x}_2\hat{x}_1} & \sigma_{\hat{x}_2}^2 \end{pmatrix} = \begin{pmatrix} 1 & \varepsilon \\ 0 & 1 \end{pmatrix} \begin{pmatrix} \sigma_{\hat{x}_{1|2}}^2 & 0 \\ 0 & \sigma_{\hat{x}_2}^2 \end{pmatrix} \begin{pmatrix} 1 & \varepsilon \\ 0 & 1 \end{pmatrix} \text{ with } |\varepsilon| \leq \frac{1}{2}.$$

It is the inverse of the unique *LDU*-decomposition of the inverse of the variance-covariance matrix of \hat{x}_1 and \hat{x}_2 . This result illustrates once again that Z_1^* tries to diagonalize the variance-covariance matrix, by bounding its triangular factor.

In order to generalize (43) to dimensions higher than two, we start from the inverse of the *LDU*-decomposition of $Q_{\hat{x}}^{-1}$. It reads $Q_{\hat{x}}^{-1} = U^{-1}D^{-1}L^{-1}$. Note that like L , L^{-1} is lower triangular having one's on the main diagonal. Also note that if the elements of L would be integer, then so would be the elements of L^{-1} . In fact, matrix L , being integer and volume-preserving, would then be the perfect candidate to diagonalize $Q_{\hat{x}}$. One would then be able to truly diagonalize $Q_{\hat{x}}$ in just one step. This observation suggests, even though the entries of L will be non-integer in general, that we choose the n -by- n matrix Z_1 as a lower triangular matrix, with integer entries and with one's on the main diagonal. Now, in order to make $Z_1^*Q_{\hat{x}}Z_1$ approximately diagonal, one could in first instance think, in analogy with the two-dimensional case, of setting Z_1 equal to L after all its elements have been rounded to their nearest integer. Unfortunately, this approach fails for the higher-dimensional case. It can namely not guarantee that *all* the non-diagonal entries of $L^{-1}Z_1$ get sufficiently bounded. Fortunately, one can do better than this by means of sweeping integer-multiples of the rows of L^{-1} . Matrix Z_1^* can then be constructed from subtracting suitable integer

multiples of the last $(n-i)$ rows of L^{-1} from row i of L^{-1} , for $i=1, \dots, (n-1)$. Using this matrix Z_1^* , one obtains

$$(44) \quad Z_1^*Q_{\hat{x}}Z_1 = (Z_1^*U^{-1})D^{-1}(L^{-1}Z_1),$$

in which, in analogy with (43), the absolute values of all non-diagonal elements of $L^{-1}Z_1$ are guaranteed to be bounded by a half. This implies, when the non-diagonal elements of L^{-1} are larger than a half in absolute value, that the diagonal entries of $Z_1^*Q_{\hat{x}}Z_1$ will be smaller than those of $Q_{\hat{x}}$. Hence, the decorrelation number will undergo an improvement through Z_1^* . But note, since (44) is the inverse of the *unique LDU*-decomposition of the inverse of $Z_1^*Q_{\hat{x}}Z_1$, that all the conditional variances stay invariant under the transformation Z_1^* . For the GPS-ambiguities this implies, not only that the variance σ_n^2 remains large, but also that the discontinuity in the spectrum of conditional variances, which is so distinctive of GPS carrier-phase processing, stays untouched. Hence, one should not expect too much from the *single* transformation Z_1^* .

6.5 Flattening the spectrum of conditional variances

In the two-dimensional case, the two type of transformations Z_1^* and Z_2^* of (37) are used in an alternating fashion in their construction of Z^* . Instead of Z_1^* and Z_2^* , one could also say that only the type Z_1^* is used, but then each time followed by a permutation of the two elements in the vector to be transformed. This suggests for the n -dimensional case, that we perform a re-ordering of the n -elements after each time that the transformation Z_1^* of the previous subsection is applied. The difficulty that we are faced with is however, what type of re-ordering to choose? That is, in dimensions higher than two, different reordering schemes are possible, all of which reduce to a simple interchange when applied to the two-dimensional case. Hence, no unambiguous generalization of the two-dimensional case seems to exist.

Fortunately, in case of GPS carrier-phase processing, already the two-dimensional scheme based on a pairwise re-ordering, allows us to obtain results that show a dramatic improvement over the original ambiguities. In order to make this clear, we first need to consider the spectrum of conditional variances. For the GPS single baseline model, based on carrier-phase data only, we have

$$(45) \quad \sigma_{\hat{x}_{j|1,\dots,n}}^2 \ll \sigma_{\hat{x}_{i|1,\dots,n}}^2 \text{ for } j=1,\dots,n-3; \quad i=n-2,n-1,n.$$

And it is this large discontinuity in the spectrum of

conditional variances, that forms a hindrance for the efficient search for the integer least-squares estimates. Our aim in constructing transformation Z^* should therefore at least be, to remove the discontinuity from the spectrum of conditional variances. And a very important consequence of such a flattening of the spectrum is, that when $n > 3$, the three large variances in the spectrum get reduced by a very significant amount. The volume-preserving property of Z^* implies namely, that the product of conditional variances remains unaffected by the transformation. Hence, by flattening the spectrum, the presence of the very small conditional variances automatically implies, that the three large variances in the spectrum have to get much smaller.

This observation now also stipulates the significance of *satellite redundancy* and *dual frequency* data. When both are absent, we have $n=3$. In that case, the absence of very small conditional variances prohibits us from "pulling down" the large variances in the spectrum. In case of satellite redundancy and/or dual frequency data however, we have $n > 3$. Now the presence of the very small conditional variances does allow us to bring the large variances in the spectrum down to much smaller values. And the larger $n-3$ is, the more we are able to bring the flattened spectrum to a lower level.

Thus, in case of GPS carrier-phase processing, a dramatic improvement can be realized, if we would be able to remove the discontinuity and enforce the spectrum of conditional variances to become much flatter. We know from subsection 6.3, see (41), that this is precisely what the transformation Z^* does for the two-dimensional case. But this suggests for the n -dimensional case, that a steep decrease in value between two consecutive conditional variances, $\sigma_{\hat{x}_{i-1}|1,\dots,n}^2$ and $\sigma_{\hat{x}_i|1,\dots,n}^2$, can be removed when the two-dimensional transformation Z^* is applied to the $(i-1)$ th and i th least-squares estimates both of which are conditioned on the last $(n-i)$ estimates: $\hat{x}_{i-1|i-1,\dots,n}$ and $\hat{x}_i|i-1,\dots,n$. Thus, instead of applying the two-dimensional transformation of subsection 6.3 to the unconditional least-squares estimates, the idea is to apply it to the *conditional* least-squares estimates. This would then give, in analogy of (41),

$$(46) \quad \sigma_{\hat{x}_{i-1}|1,\dots,n}^2 \geq \frac{3}{4} \sigma_{\hat{x}_i|1,\dots,n}^2.$$

Note that the other conditional variances remain unaffected by the transformation. This is simply a consequence of the fact that we are transforming *conditional* least-squares estimates. Result (46) implies for $i=n-2$, that we are able to

close the large gap between the $(n-3)$ th and $(n-2)$ th conditional variance. Ofcourse, after the transformation has been applied, other, but smaller, discontinuities emerge. For instance, if the transformation has been applied for $i=n-2$, then $\sigma_{\hat{x}_{n-2}|1,\dots,n}^2 < \sigma_{\hat{x}_{n-1}|1,\dots,n}^2$ and $\sigma_{\hat{x}_{n-1}|1,\dots,n}^2 < \sigma_{\hat{x}_n|1,\dots,n}^2$. But, also they can be removed by applying the two-dimensional transformation. In fact, one can continue in this way and flatten the complete spectrum of conditional variances.

In summary, the proposed method thus flattens the n -dimensional spectrum of conditional variances through a repeated application of the two-dimensional transformation Z^* to the conditional least-squares estimates. And the larger $n-3$ is, the lower the level of the transformed spectrum. As a result the n th ambiguity shows a dramatic improvement in precision, $\sigma_{\hat{x}_n}^2 \ll \sigma_{\hat{x}_n}^2$. And this can also be assured for the remaining ambiguities, when the low level of the transformed spectrum is combined with a bounding of the triangular factor, as given in the previous subsection. Thus the proposed method of ambiguity-transformation, returns less correlated ambiguities with a significantly improved precision and allows for a very efficient search for the transformed integer least-squares ambiguities. Our numerical experiments indicate for instance, that with dual-frequency carrier-phase data, based on two epochs of data, with a one second sampling interval, standard deviations of the transformed ambiguities are obtained that usually are well below the one cycle level.

7 Summary and concluding remarks

In this contribution a new method was introduced for computing the integer least-squares estimates of the GPS ambiguities. It was shown in section 3 that this problem can be reduced to the integer least-squares problem

$$(47) \quad \min_x (\hat{x} - x)^T Q_x^{-1} (\hat{x} - x), \quad x \in Z^n.$$

In case of GPS however, when short timespan carrier-phase data is used, the elements of \hat{x} , being the least-squares estimates of the double-difference ambiguities, are extremely correlated. Also the confidence ellipsoid will then be extremely elongated. It is not uncommon for instance, to have an elongation in the order of 10^4 , when the data is based on two epochs, one second apart (if length of minor axis 1 cm, then length of major axis 100 mtr.). The amount of correlation between the ambiguities and therefore the non-diagonality of Q_x can be measured by the determinant r_x^2

of the correlation matrix. When $r_{\hat{x}}$ equals one, $Q_{\hat{x}}$ is diagonal, and when $r_{\hat{x}}$ is close to zero, then $Q_{\hat{x}}$ is far from diagonal. The scalar $r_{\hat{x}}^2$ reads in terms of the conditional variances and unconditional variances, as

$$(48) \quad r_{\hat{x}}^2 = \prod_{i=1}^n \frac{\sigma_{\hat{x}_{i+1:n}}^2}{\sigma_{\hat{x}_i}^2}.$$

And because of the large discontinuity in the spectrum of conditional variances of the ambiguities, $r_{\hat{x}}$ is usually extremely small. For instance, with dual-frequency data and a satellite redundancy of only one, $r_{\hat{x}}$ can be in the order of 10^{-19} , when the data is based on two epochs of data, one second apart. As a result of this extreme non-diagonality of $Q_{\hat{x}}$, the efficiency in solving the above integer least-squares problem is severely hindered. The idea is therefore to reparametrize the above integer least-squares problem, such that an equivalent formulation is obtained, but one that is much easier to solve. By introducing the reparametrization

$$(49) \quad \hat{z} - z = Z^*(\hat{x} - x),$$

with Z^* being integer and volume-preserving, we obtain the equivalent minimization problem

$$(50) \quad \min_z (\hat{z} - z)^* Q_{\hat{z}}^{-1} (\hat{z} - z), \quad z \in Z^n,$$

with the new variance-covariance matrix $Q_{\hat{z}} = Z^* Q_{\hat{x}} Z$. The transformed integer least-squares problem becomes trivial, when the new variance-covariance matrix $Q_{\hat{z}}$ is diagonal. That is, when $r_{\hat{z}}$ equals one. The idea is therefore to come up with a matrix Z^* that allows $r_{\hat{z}}$ to be close to one. Based on integer approximating the conditional least-squares transformation, the construction of such a matrix was given in subsection 6.3 for the two-dimensional case. It returns ambiguities with an improved precision and guarantees, because of $r_{\hat{z}}^2 \geq \frac{3}{4}$, that

$$(51) \quad \sigma_{\hat{z}_i}^2 \geq \frac{3}{4} \sigma_{\hat{z}_i}^2.$$

From this it followed, that one can remove the discontinuity in the spectrum of conditional variances of the original ambiguities, by means of a repeated application of the two-dimensional transformation to the sequential conditional least-squares estimates of the ambiguities. As a result, the method returns: less correlated ambiguities (for instance an improvement from the above given $r_{\hat{x}} \cong 10^{-19}$ to $r_{\hat{z}} \cong 0.5$); significantly more precise ambiguities (standard deviations of the transformed ambiguities that are well below the one cycle level are not uncommon, when dual-frequency carrier-phase data is used, based on two epochs of data, with a one second

interval); and it allows one to perform the search for the transformed integer least-squares ambiguities, based on

$$(52) \quad (z_i - \hat{z}_{i|1:n})^2 \leq \lambda_{\hat{z}_i} \sigma_{\hat{z}_{i+1:n}}^2 \chi^2,$$

in a highly efficient manner.

To conclude, we finish with a few remarks on some untouched GPS issues. As it was pointed out in the introduction, the proposed method is directed towards solving the estimation problem of GPS-ambiguity fixing. But still, it also bears some relation to the validation step. This is particular true, when the validation step is based on a comparison of the most-likely and the second most-likely integer minimizer. With some minor changes in the search, the method namely also allows an efficient computation of the second most-likely integer minimizer. Also, since the method is significantly faster than existing methods for GPS-ambiguity fixing, the gain in efficiency leaves us room for handling higher dimensional integer least-squares problems. This may in particular be useful in a GPS-network approach. And finally, we would like to point out that the given ambiguity transformation is completely determined by the variance-covariance matrix of the ambiguities. Even the a posteriori variance-factor need not be known. This stipulates that actual measurements are not needed to perform the transformation to near-diagonality. Hence, the necessary computations can be done in principle at the designing stage, prior the actual measurement stage.

8 References

1. Baarda, W. (1967): *Statistical Concepts in Geodesy*, Netherlands Geodetic Commission, Publication on Geodesy, New Series, Vol. 2, No. 4.
2. Blewitt, G. (1989): Carrier Phase Ambiguity Resolution for the Global Positioning System Applied to Geodetic Baselines up to 2000 km. *Journal of Geophysical Research*, Vol. 94, No. B8, pp. 10.187-10.203.
3. Cocard, M. and A. Geiger (1992): Systematic Search for all Possible Widelines. Paper presented at the *Sixth International Geodetic Symposium on Satellite Positioning*, Columbus, Ohio, March 17-20, 1992.
4. Counselman, C.C. and S.A. Gourevitch (1981): Miniature Interferometer Terminals for Earth Surveying: Ambiguity and Multipath with Global Positioning System. *IEEE Transactions on Geoscience and Remote Sensing*, Vol. GE-19, No. 4, pp. 244-252.
5. Euler, H.J. and H. Landau (1992): Fast GPS ambiguity

- resolution on-the-fly for real-time applications. Paper presented at the *Sixth International Geodetic Symposium on Satellite Positioning*, Columbus, Ohio, March 17-20, 1992.
6. Frei, E. (1991): *Rapid Differential Positioning with the Global Positioning System*. Schweizerischen Geodätischen Kommission, Band 44.
 7. Hatch, R. (1991): Instantaneous Ambiguity Resolution. Proceedings of IAG *International Symposium 107 on Kinematic Systems in Geodesy, Surveying and Remote Sensing*, Sept. 10-13, 1990, Springer Verlag, New York, pp. 299-308.
 8. Kleusberg, A. (1990): A Review of Kinematic and Static GPS Surveying Procedures. Proceedings of the *Second International Symposium on Precise Positioning with the Global Positioning System*, Ottawa, Canada, September 3-7 1990, pp. 1102-1113.
 9. Lenstra, H.W. (1981): *Integer Programming with a Fixed Number of Variables*. University of Amsterdam, Dept. of Mathematics, report 81-03.
 10. Remondi, B.W. (1986): Performing Centimeter-Level Surveys in Seconds with GPS Carrier Phase: Initial Results. *Journal of Navigation*, Volume III, The Institute of Navigation.
 11. Scheffé, H. (1956): *The Analysis of Variance*. John Wiley and Sons.
 12. Tienstra, J.M. (1956): *Theory of the Adjustment of Normally Distributed Observations*, Argus, Amsterdam.
 13. Teunissen, P.J.G. (1990): GPS op afstand bekeken. In: *Een halve eeuw in de goede richting*. Lustrumboek Snellius 1985-1990, pp. 215-233.
 14. Teunissen, P.J.G. (1990): Nonlinear Least Squares. *Manuscripta Geodaetica*, 15, pp. 137-150.
 15. Teunissen, P.J.G. (1993): Integer Least-Squares: Fast Estimation of the Carrier Phase Ambiguities. *First International Symposium on the Mathematical and Physical Foundations of Geodesy*, September 7-9 1993, Invited paper (in print), Stuttgart.
 16. Teunissen, P.J.G. (1993): *The Invertible GPS Ambiguity Transformations*, Delft Geodetic Computing Centre (LGR), 9 p.
 17. Wübbena, G. (1991): *Zur Modellierung van GPS-Beobachtungen für die hochgenaue Positionsbestimmung*, Hannover, 1991.