

Detector Alignment by Tracks with Millepede II

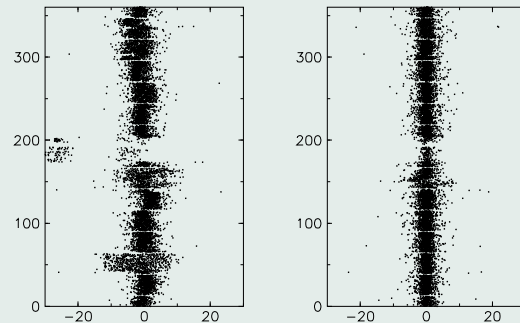
Volker Blobel – Universität Hamburg

Abstract

Tracking detectors in high energy physics experiments require an accurate determination of a large number of alignment parameters in order to allow a precise reconstruction of tracks and vertices. In addition to the initial optical survey and corrections for electronics and mechanical effects the use of tracks in a special software alignment is essential. The general program MILLEPEDE performs a simultaneous least squares fit of a large number of tracks, with determination of (global) alignment parameters and (local) track parameters. The version II, which is under development, provides several options for the solution of large matrix equations. Aim is a program that is able to determine up to 100 000 alignment parameters in a reasonable time on a standard PC.

1. Introduction
2. Alignment of a toy detector
3. Millepede I (1997 -)
4. Millepede II (2005 -)

Summary



1. Introduction

Purpose of instrument calibration: Instrument calibration is intended to **eliminate or reduce bias** in an instrument's readings over a range for all continuous values. For this purpose, **reference standards** with known values for selected points covering the range of interest are measured with the instrument in question. Then a **functional relationship** is established between the values of the standards and the corresponding measurements [...from NIST].

Alignment/calibration of HEP track detectors:

- based **only** on track residual minimization (*incomplete* data, several degrees of freedom undefined) and on survey data;
- no “reference standards with known values” exist (\approx exceptions are data from $e^+e^- \rightarrow \mu^+\mu^-$ and cosmics without field).

Alignment/calibration requires to *understand* the detector (functional relationship) and to optimize thousands or ten thousands of parameters. Aim is – after an initial optical survey and corrections for electronics and mechanical effects:

- reduce χ^2 of the track fits, in order to improve track and vertex recognition, and
- increase precision of reconstructed tracks and vertices, eliminating or reducing bias in detector data.

... essential for important aspects of physics analysis with large accurate vertex detectors with potential precision of a few μm .

Optimization

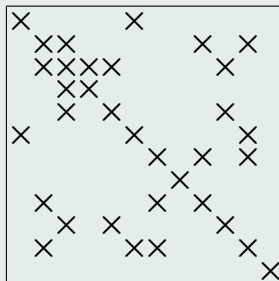
Alignment as an optimization problem: The standard method for the determination of a large number of *correlated* parameters results is the definition of an objective function $F(\Delta\mathbf{p})$, and the minimization of the function by the solution of a large system of linear equations ($\Delta\mathbf{p}$ = corrections to alignment parameters)

$$\boxed{\mathbf{C} \Delta\mathbf{p} = \mathbf{b}}$$

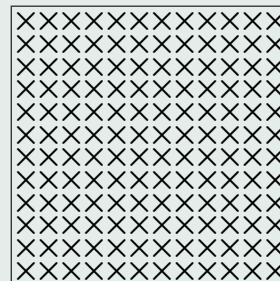
\mathbf{C} is a sparse symmetric $n \times n$ matrix, with $n = \dots 100\,000$

- Construction of correct matrix \mathbf{C} and vector \mathbf{b} from the data, or at least of a good approximation to the matrix \mathbf{C} and vector \mathbf{b} , and ...
- correct solution of the large matrix equation, or at least a good approximation to the correct solution.

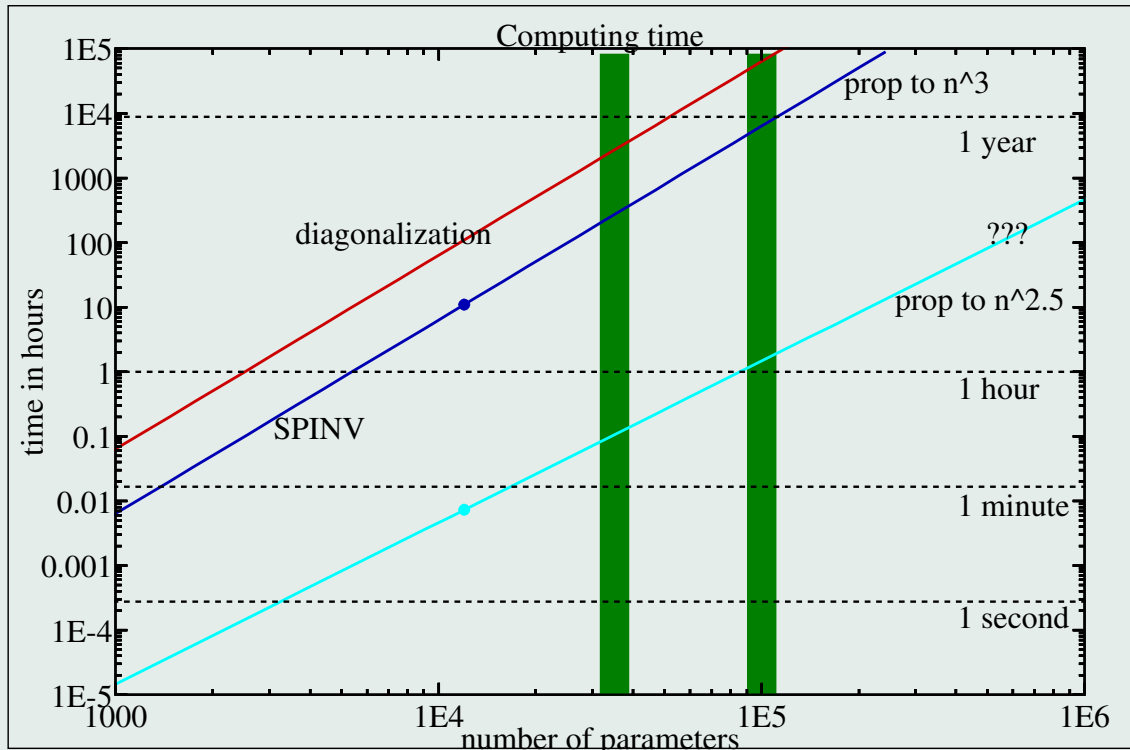
The inverse of a sparse matrix



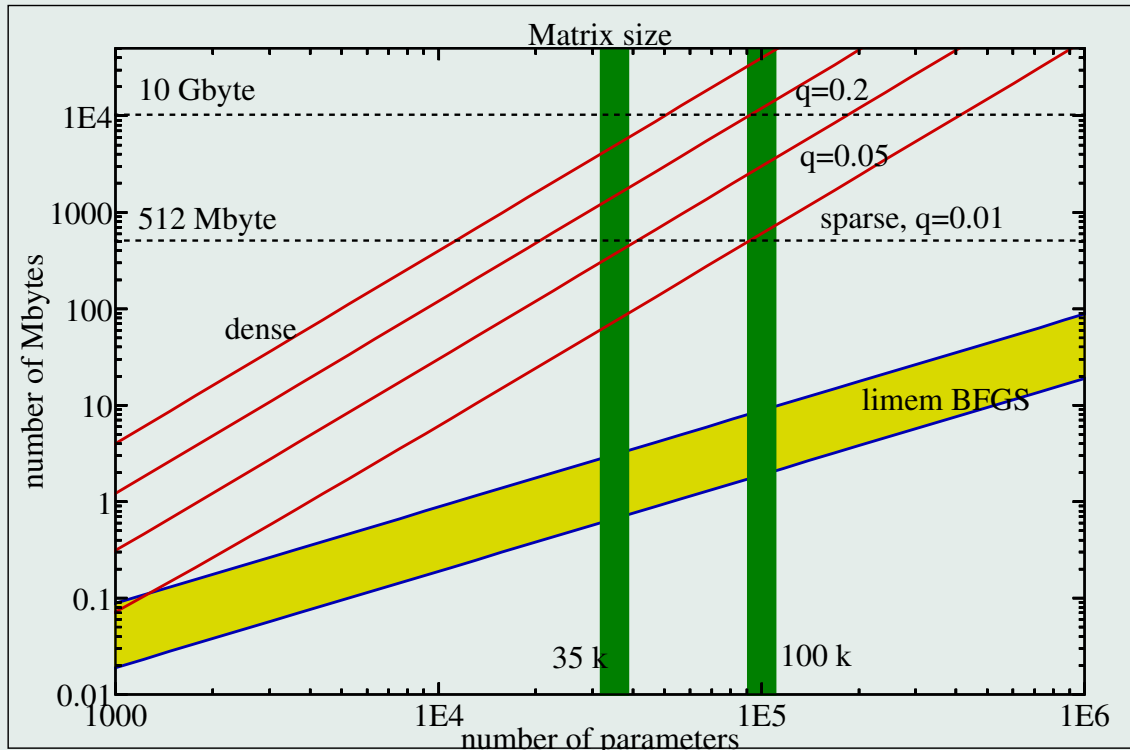
is a dense matrix:



Calculation of the – dense – matrix \mathbf{C}^{-1} requires $(n^2 + n)/2$ words and a time $\propto n^3 \Rightarrow$ solution by inversion is (almost) impossible \Rightarrow use other methods.



The two circles are points, measured with Millepede II for $n = 12\,015$. The lower curve is for an iterative method, which requires only space for the sparse matrix C .



q = fraction of non-zero off-diagonal elements

Note: printing the elements of a n -by- n matrix for 100 000 requires $\approx 10 \text{ m}^3$ paper (double-sided printing).

Direct solution: no iterations necessary (but inversion time $\propto n^3$)

Rate of convergence of iterative methods: $\{\mathbf{x}_k\}$ be sequence in \mathbb{R}^n that converges to solution \mathbf{x}^*

$$\begin{array}{ll} \text{linear convergence} & \frac{\|\mathbf{x}_{k+1} - \mathbf{x}^*\|}{\|\mathbf{x}_k - \mathbf{x}^*\|} \leq r \quad r \in (0, 1) \\ \text{superlinear convergence} & \lim_{k \rightarrow \infty} \frac{\|\mathbf{x}_{k+1} - \mathbf{x}^*\|}{\|\mathbf{x}_k - \mathbf{x}^*\|} = 0 \\ \text{quadratic convergence} & \frac{\|\mathbf{x}_{k+1} - \mathbf{x}^*\|}{\|\mathbf{x}_k - \mathbf{x}^*\|^2} \leq M \quad M \text{ positive} \end{array}$$

for all k sufficiently large.

Speed of convergence depends on the eigenvalue spectrum – slow convergence for small eigenvalues.

Note: an iterative method can be faster than a direct method!

Rate of convergence

Linear convergence: Example is steepest-descent;

constant r may be close to 1; iteration numbers of 100 or 1000 not uncommon; an algorithm with r close to 1, i.e. $1 - r \ll 1$ is in practice considered as not converging at all. For a quadratic objective function:

$$\frac{\|\mathbf{x}_{k+1} - \mathbf{x}^*\|}{\|\mathbf{x}_k - \mathbf{x}^*\|} \leq r = \frac{\kappa - 1}{\kappa + 1} \quad \kappa = \text{Condition} = \frac{\lambda_{\max}}{\lambda_{\min}} \quad (\gg 1)$$

Often in applications a cut on small $\|\mathbf{x}_{k+1} - \mathbf{x}_k\|$ is used, which means

$$\|\mathbf{x}_{k+1} - \mathbf{x}_k\| \approx (1 - r)\|\mathbf{x}_k - \mathbf{x}^*\| \quad \text{with } (1 - r) \ll 1$$

Superlinear convergence: example is Quasi-NEWTON method;

Quadratic convergence: example is NEWTONs method;

roughly the number of exact digits double at each iteration; typically 3 to 5 iterations needed.

Least squares

Minimize objective function $F(\mathbf{q}) = \frac{1}{2} \sum_{j=1}^m r_j^2(\mathbf{q})$

The derivatives of $F(\mathbf{q})$ can be expressed in terms of the Jacobian of \mathbf{r} (m -by- n matrix of first partial derivatives) defined by

$$\mathbf{J}(\mathbf{q}) = \left[\frac{\partial \mathbf{r}_j}{\partial q_i} \right]_{j=1, \dots, m; i=1, \dots, n}$$

Exploit the special structure of least squares:

$$\nabla F(\mathbf{q}) = \mathbf{J}(\mathbf{q})^T \mathbf{r}(\mathbf{q}) \quad \nabla^2 F(\mathbf{q}) = \mathbf{J}(\mathbf{q})^T \mathbf{J}(\mathbf{q}) \left[+ \sum_{j=1}^m r_j(\mathbf{q}) \nabla^2 r_j(\mathbf{q}) \right]$$

where the contribution in brackets [...] is usually neglected (Gauss-Newton method). Note that the full matrix \mathbf{J} has never to be stored: one row at a time is sufficient.

For a linear problem (fixed first partial derivatives) only one step is necessary to determine \mathbf{q} :

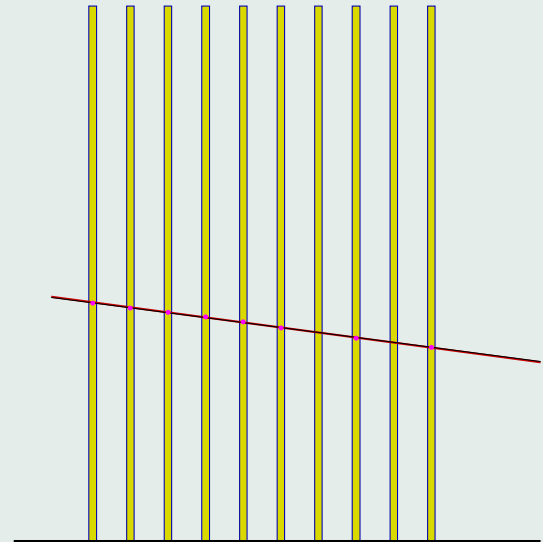
$$\underbrace{(\mathbf{J}(\mathbf{q})^T \mathbf{J}(\mathbf{q}))}_{\text{symmetric matrix}} \mathbf{q} = \underbrace{(\mathbf{J}(\mathbf{q})^T \mathbf{r}(\mathbf{q}))}_{\text{vector}}$$

Alignment problems are usually almost linear \Rightarrow one least squares step may be sufficient ...
... but there may be other reasons for iterations.

2. Alignment of a toy detector

Test of alignment method with a MC toy track detector model:

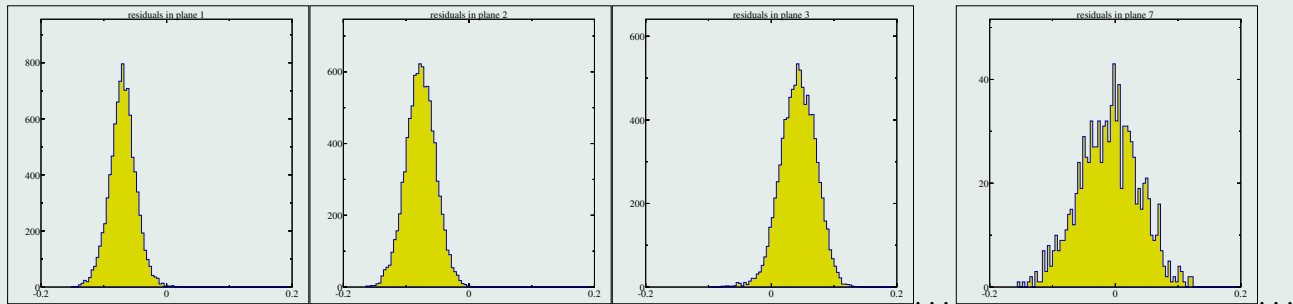
- 10 planes of tracking chambers, 1 m high, 10 cm distance, no magnetic field;
- accuracy $\sigma \approx 200\mu\text{m}$, with efficiency $\epsilon = 90\%$;
- plane 7 sick: accuracy $\sigma \approx 400\mu\text{m}$, with efficiency $\epsilon = 10\%$;
- 10 000 tracks with 82 000 hits available for alignment;
- **Misalignment:** the vertical position of the chambers are displaced by $\approx 0.1\text{cm}$ (normal distributed).



First attempt based on residuals

The first alignment attempt is based on the distribution of hit residuals:

- A straight line is fitted to the track data.
- The residuals (= measured vertical coordinate — fitted coordinate) are histogrammed, separately for each plane.



- The mean value of the residuals is taken as correction to the vertical plane position.

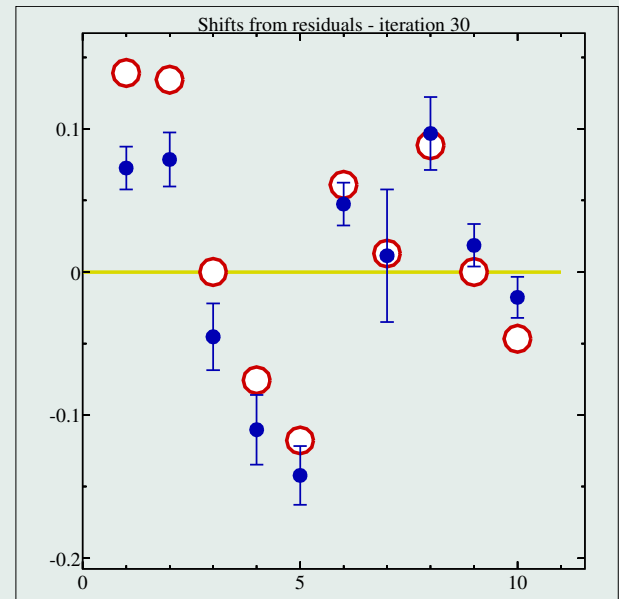
This is the standard method used in many experiments. What is the convergence behaviour?

Result from the first attempt

Large changes in first iteration, small changes in second iteration, almost no progress afterwards.

After 30 iterations ...

ID	true shift	determined	mean residual
1	0.1391	0.0727	0 ± 150
2	0.1345	0.0786	0 ± 189
3	0.0000	-0.0453	0 ± 234
4	-0.0756	-0.1102	0 ± 244
5	-0.1177	-0.1422	0 ± 205
6	0.0610	0.0475	0 ± 150
7	0.0130	0.0114	0 ± 464
8	0.0886	0.0968	0 ± 255
9	0.0000	0.0186	0 ± 149
10	-0.0467	-0.0176	0 ± 143



red circle = true shift (displacement)

blue disc = displacement, determined from residuum

First attempt – Discussion

The result is not (yet) encouraging!

The reason for non-convergence is simple:

Two degrees of freedom are undefined: a simultaneous shift and a overall shearing of the planes!

(... this simple fact is not always mentioned in reports on the method.)

Improvement for second residual attempt:

Fix the displacement (i.e. $\text{displacement} = 0$) of two planes, which are assumed to be carefully aligned externally (e.g. planes 3 and 9).

Other possibilities are:

- Use only fixed planes (planes 3 and 9) in the fit, and determine the residuals of other planes;
- for the determination of the displacement of a certain plane use all other planes in the fit.

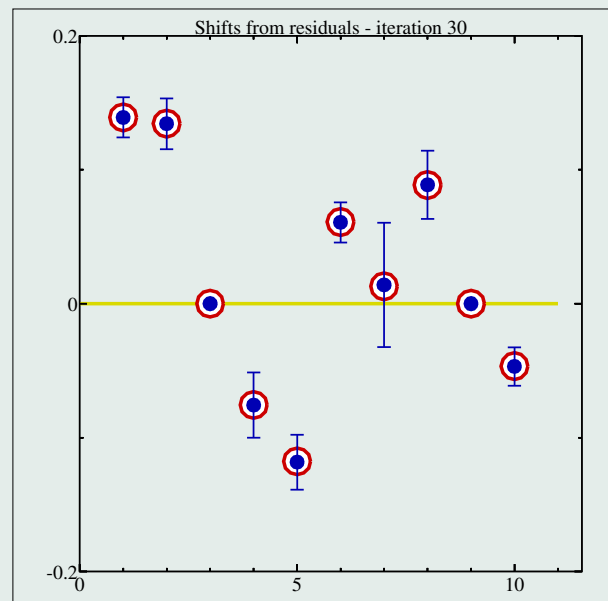
These possibilities are in fact used by several collaborations!

Results from the second attempt

Large changes in first iteration, then many smaller and smaller changes: convergence is **linear** and slow, because the determination of displacements is based on biased fits.

After 30 iterations with planes 3 and 9 fixed (displacement = 0) ...

ID	true shift	determined	mean residual
1	0.1391	0.1391	-1 ± 150
2	0.1345	0.1344	0 ± 189
3	0.0000	0	2 ± 234
4	-0.0756	-0.0758	0 ± 244
5	-0.1177	-0.1183	0 ± 205
6	0.0610	0.0607	0 ± 150
7	0.0130	0.0140	0 ± 464
8	0.0886	0.0888	0 ± 255
9	0.0000	0	0 ± 149
10	-0.0467	-0.0469	0 ± 143



red circle = true shift (displacement)

blue disc = displacement, determined from residuum

Use of higher mathematics?

Residual-based methods work with biased results. Can the bias be avoided by an improved fit?

Yes: include the alignment parameters in the parameters fitted in track fits – requires a simultaneous fits of many tracks, with determination of (global) alignment parameters and (local) track parameters.

model: $y_i \cong a_1^{\text{local}} + a_2^{\text{local}} \cdot x_i + a_j^{\text{global}}$ a_j^{global} = shift for plane j , where y_i is measured

1 tracks	2 + 10 = 12 parameters	9 equations
2 tracks	4 + 10 = 14 parameters	18 equations
...
10 000 tracks	20 010 parameters	82 000 equations

... a linear least squares problem of $m = 82\,000$ equations (measurements) and $n = 20\,010$ parameters with $n \ll m$, which requires the solution of a matrix equation with 20010-by-20010 matrix.

The MILLEPEDE principle allows to reduce the least squares problem to $n = 10$ parameters!

Results from a simultaneous fit

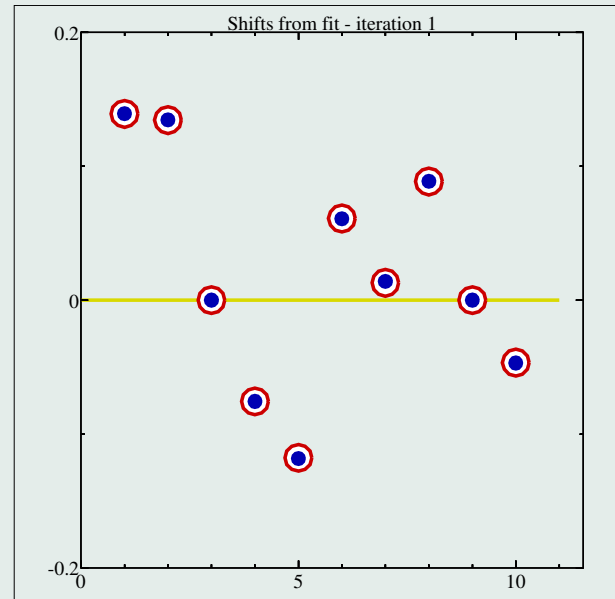
After one step (with planes 3 and 9 fixed at displacement = 0) ...

ID	true shift	determined	ρ	mean residual
1	0.1391	0.1393 ± 0.0004	0.68	0 ± 150
2	0.1345	0.1346 ± 0.0003	0.66	0 ± 189
3	0.0000			0 ± 234
4	-0.0756	-0.0756 ± 0.0003	0.58	0 ± 244
5	-0.1177	-0.1182 ± 0.0003	0.53	0 ± 205
6	0.0610	0.0608 ± 0.0003	0.50	0 ± 150
7	0.0130	0.0141 ± 0.0007	0.20	0 ± 464
8	0.0886	0.0888 ± 0.0003	0.53	0 ± 255
9	0.0000			0 ± 149
10	-0.0467	-0.0469 ± 0.0003	0.57	0 ± 143

(ρ = global correlation coefficient)

red circle = true shift (displacement)

blue disc = displacement, determined in fit



One step is sufficient: 1. step $\Delta\chi^2 = 1.277 \times 10^6$

2. step $\Delta\chi^2 = 1.159 \times 10^{-5}$

Determination of drift velocities

... 10 additional parameters

Improvement: include, in addition, corrections to the drift velocities for each plane: $\Delta v_{\text{drift}}/v_{\text{drift}}$

$$y_i \cong a_1^{\text{local}} + a_2^{\text{local}} \cdot x_i + a_j^{\text{global}} + \ell_{\text{drift},i} \cdot \left(\frac{\Delta v_{\text{drift}}}{v_{\text{drift}}} \right)_j$$

a_j^{global} = shift for plane j

$\left(\frac{\Delta v_{\text{drift}}}{v_{\text{drift}}} \right)_j$ = relative v_{drift} difference

reduction of residual σ by 30 - 40 %

ID	true shift	determined	ρ	$\Delta v_{\text{drift}}/v_{\text{drift}}$	determined	ρ	mean residual
1	0.1391	0.1393 ± 0.0004	0.68	0.0020	0.0019 ± 0.0002	0.016	0 ± 119
2	0.1345	0.1346 ± 0.0003	0.66	-0.0153	-0.0150 ± 0.0002	0.020	0 ± 128
3	0.0000			0.0193	0.0194 ± 0.0002	0.017	0 ± 137
4	-0.0756	-0.0756 ± 0.0003	0.58	0.0200	0.0197 ± 0.0002	0.013	0 ± 139
5	-0.1177	-0.1182 ± 0.0003	0.53	-0.0138	-0.0136 ± 0.0002	0.013	0 ± 141
6	0.0610	0.0608 ± 0.0003	0.50	0.0003	0.0004 ± 0.0002	0.019	0 ± 139
7	0.0130	0.0141 ± 0.0007	0.20	-0.0306	-0.0303 ± 0.0006	0.038	0 ± 348
8	0.0886	0.0888 ± 0.0003	0.53	0.0237	0.0238 ± 0.0002	0.018	0 ± 134
9	0.0000			-0.0044	-0.0044 ± 0.0002	0.008	0 ± 127
10	-0.0467	-0.0469 ± 0.0003	0.57	0.0021	0.0019 ± 0.0002	0.013	0 ± 117

... this would be rather difficult with a pure residual-based method.

The next improvement would be the introduction of wire T_0 's – additional $10 \times 25 = 250$ parameters.

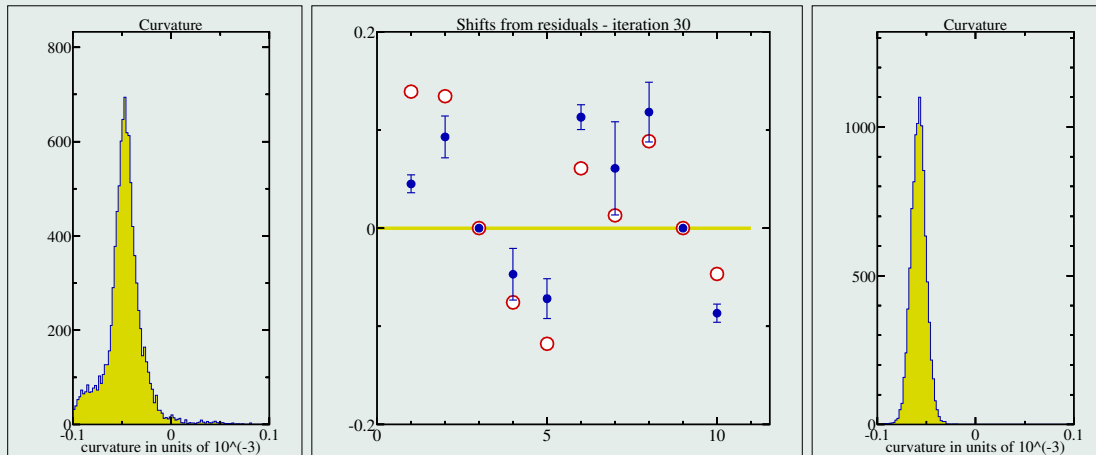
A more realistic scenario

So far the tracks are fitted with a straight line. Now a third parameter is added to the parameterization and a parabola is fitted (i.e. case of unknown momentum):

$$\text{model:} \quad y_i \cong a_1^{\text{local}} + a_2^{\text{local}} \cdot x_i + a_3^{\text{local}} \cdot x_i^2$$

The initial misalignment will create a curvature $\neq 0$ – the data are now insufficient to determine the true shifts.

Alignment by residual method



fitted curvature before ...

... after alignment

3. Millepede I

Determination of corrections $\Delta \mathbf{p}$ for alignment parameters \mathbf{p} is based on minimization of residuals – the difference between fitted and measured track position:

$$\Delta_i = \text{fitted value} - \text{measured value}$$

A “global” objective function $F(\Delta \mathbf{p}, \mathbf{q})$ or χ^2 -function is constructed, which depends on the corrections $\Delta \mathbf{p}$ and **all** track parameters \mathbf{q}

$$F(\Delta \mathbf{p}, \mathbf{q}) \equiv \chi^2(\Delta \mathbf{p}, \mathbf{q}_j) = \sum_{\text{data sets}} \left(\sum_{\text{events}} \left(\sum_{\text{tracks}} \left(\sum_{\text{hits}} \Delta_i^2 / \sigma_i^2 \right) \right) \right)$$

Data sets are

- Physics data: track data from e^+e^- , e^-p , pp -reactions,
- Cosmics with magnetic field (large distance to IP) and without magnetic field (straight tracks, curvature zero),
- vertex- or mass-constrained track data,
- external alignment data.

A mixture of different data is recommended, in order to introduce different correlations between the alignment parameters and to increase the precision.

Minimization requires first derivatives of hit residuals w.r.t. the parameters – *linearization*. Residuals and derivatives are used to determine matrix and vector of the so-called normal equations of least squares:

$$\left(\begin{array}{c|ccc} \mathbf{C} & \dots & \mathbf{H}_k & \dots \\ \hline \dots & \dots & \mathbf{0} & \mathbf{0} \\ \mathbf{H}_k^T & \mathbf{0} & \mathbf{C}_k^{\text{track}} & \mathbf{0} \\ \dots & \mathbf{0} & \mathbf{0} & \dots \end{array} \right) \times \begin{pmatrix} \Delta \mathbf{p} \\ \dots \\ \Delta \mathbf{q}_k^{\text{track}} \\ \dots \end{pmatrix} = \begin{pmatrix} \mathbf{b} \\ \dots \\ \mathbf{b}_k^{\text{track}} \\ \dots \end{pmatrix} \quad k = \text{track index}$$

which defines the alignment corrections $\Delta \mathbf{p}$ and all track parameters \mathbf{q}_k from the data used in the fit. The size of the huge matrix depends on the number of alignment parameters **and** on number of tracks. Note: ignoring the alignment parameters, the (biased) track parameters are determined from the matrix equation

$$\mathbf{C}_k^{\text{track}} \Delta \mathbf{q}_k^{\text{track}} = \mathbf{b}_k^{\text{track}} \quad \text{with cov. matrix} \quad \mathbf{V}_k^{\text{track}} = (\mathbf{C}_k^{\text{track}})^{-1}$$

with a small matrix and vector – or ignore track fits and calculate

$$\mathbf{C} \Delta \mathbf{p} = \mathbf{b}$$

The Millepede principle allows to find the solution of the huge equation, because the sparse structure of the matrix allows a reduction to the matrix size for the alignment parameters, using the results from the track fits and the rectangular *mixed* matrices \mathbf{H}_k .

Normal equations

Simultaneous least squares fit of all global and all local parameters (i.e. all tracks).

$$k'\text{th track:} \quad y_i \cong f(x_i; \mathbf{p}^{\text{global}}, \mathbf{q}^{\text{local}}) + \left(\mathbf{d}_i^{\text{global}} \right)^T \Delta \mathbf{p}^{\text{global}} + \left(\boldsymbol{\delta}_i^{\text{local}} \right)^T \Delta \mathbf{q}_k^{\text{local}}$$

The complete matrix equation for global and local parameters includes sums over track index k and contains many matrices: n -by- n matrices \mathbf{C} for n global parameters and m -by- m matrices $\mathbf{C}_k^{\text{local}}$ and n -by- m matrices $\mathbf{H}_k^{\text{global-local}}$

$$\begin{pmatrix} \sum_k \mathbf{C}_k^{\text{global}} & \cdots & \mathbf{H}_k^{\text{global-local}} & \cdots \\ \vdots & \ddots & 0 & 0 \\ \left(\mathbf{H}_k^{\text{global-local}} \right)^T & 0 & \mathbf{C}_k^{\text{local}} & 0 \\ \vdots & 0 & 0 & \ddots \end{pmatrix} \times \begin{pmatrix} \Delta \mathbf{p}^{\text{global}} \\ \vdots \\ \Delta \mathbf{q}_k^{\text{local}} \\ \vdots \end{pmatrix} = \begin{pmatrix} \sum_k \mathbf{b}_k^{\text{global}} \\ \vdots \\ \mathbf{b}_k^{\text{local}} \\ \vdots \end{pmatrix}$$

If the $\mathbf{H}_k^{\text{global-local}}$ are neglected, the complete equation decays into $1 + K$ independent matrix equations.

$\Delta \mathbf{p}^{\text{global}}$ can be calculated without approximation with a great simplification: \Rightarrow

For each track in a loop, on all tracks:

1. **Track- or other fit:** perform fit by finding the best local parameter values for the actual track until convergence with determination of the covariance matrix \mathbf{V}_k of the local parameters
2. **Derivatives:** calculate for all hits (index i) the vectors of derivatives δ_i^{local} and d_i^{global} for all local and the relevant global parameters, and update matrices:

$$\mathbf{C} := \mathbf{C} + \sum_i w_i d_i^{\text{global}} \left(d_i^{\text{global}} \right)^T \quad \mathbf{b} := \mathbf{b} + \sum_i w_i r_i d_i^{\text{global}} \quad \mathbf{H}_k = \sum_i w_i d_i^{\text{global}} \left(\delta_i^{\text{local}} \right)^T$$

and finally for the track $\mathbf{C} := \mathbf{C} - \mathbf{H}_k \mathbf{V}_k \mathbf{H}_k^T \quad \mathbf{b} := \mathbf{b} - \mathbf{H}_k (\mathbf{V}_k \mathbf{b}_k)$

The 'blue' equations transfer the 'local' information to the global parameters.

After the loop on all tracks the complete information is collected; now the matrix equation for the global parameters has to be solved:

$$\text{solve } \mathbf{C} \Delta \mathbf{p}^{\text{global}} = \mathbf{b} \quad \text{for } \Delta \mathbf{p}^{\text{global}} \quad \text{e.g. by } \Delta \mathbf{p}^{\text{global}} = \mathbf{C}^{-1} \mathbf{b}$$

Note: matrices \mathbf{C} and vectors \mathbf{b} from several data sets can be simply added to get combined result.

Numerical linear algebra

Determination of *good* solution for $\Delta \mathbf{p}$ in

$$\boxed{C \Delta \mathbf{p} = \mathbf{b}} \quad C^{-1} = V = \text{covariance matrix for parameters}$$

with large or huge $n \times n$ matrix depends on **algorithm** and **data**, and requires

- 0. **Scaling:** Use consistent units in data and variables – matrix and vector should be well scaled, i.e. elements should have similar precision.
- 1. **Algorithm – Stability:** With a stable algorithm the computed solution is the exact solution of a nearby problem. Gaussian Elimination with diagonal pivoting (restricted scan for largest next element) is considered to be a stable algorithm for positive definite matrices.
- 2. **Data – Conditioning:** Data are called ill-conditioned, if small changes in the data can cause large changes in the solution (\rightarrow small eigenvalue(s)). Matrix is ill-conditioned if variables are undefined or poorly defined or strongly correlated.

Standard method for solution of $\mathbf{C}\Delta\mathbf{p} = \mathbf{b}$ with symmetric matrix is stable Gauss algorithm with pivot selection of diagonal, with

$$\text{Computing time} = \text{constant} \times n^3 \quad < 1 \text{ hour for } n = \text{several thousand}$$

(in reality “constant” increases for large n).

A standard matrix routine will fail – at least a few parameters out of many thousands will be badly defined. – Matrix is destroyed without result for $\Delta\mathbf{p}$.

Subroutine SPMINV (in Millepede) for symmetric matrices in $(n^2 + n)/2$ words: choose largest pivot, but stop inversion if no acceptable pivot found, i.e. invert largest possible submatrix; return zero corrections for remaining parameters.

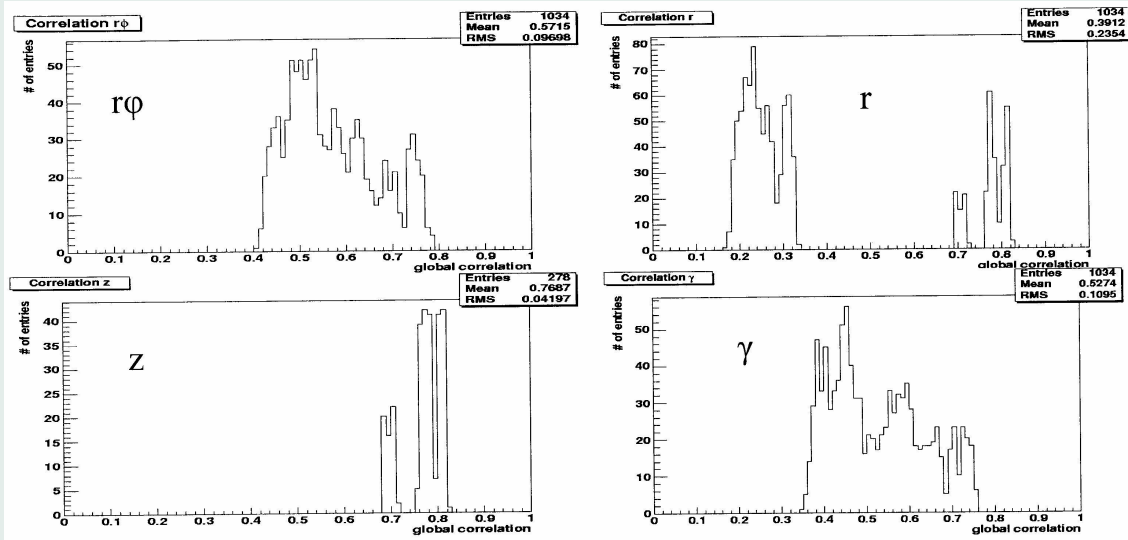
All variances and covariances available in inverse matrix.

The **global correlation coefficient**, ρ_j is a measure of the total amount of correlation between the j -th parameter and *all* the other variables. It is the largest correlation between the j -th parameter and every possible linear combination of all the other variables.

$$\rho_j = \sqrt{1 - \frac{1}{(\mathbf{V})_{jj} \cdot (\mathbf{C})_{jj}}} \quad \text{and} \quad (\mathbf{V})_{jj} \cdot (\mathbf{C})_{jj} = \frac{1}{1 - \rho_j^2} \quad \mathbf{V} = \mathbf{C}^{-1}$$

Global correlations

Range of global correlation coefficient is $0 \dots 1$.



Values ≈ 1 means strong correlation and almost singular matrix – inversion may be impossible.

Values depend on geometry and type of data – additional data (cosmics, vertex constrained tracks) can reduce the global correlation and improve the alignment.

Undefined degrees of freedom ... or weakly defined degrees of freedom

Alignment of HEP track detectors ... if based **only** on track residual minimization: *incomplete* data, with several degrees of freedom undefined! Certain parameters are **undefined** or only weakly defined and **could distort the detector**.

General **linear** transformation of whole detector with translation and 3×3 matrix **R**

$$\begin{pmatrix} x' \\ y' \\ z' \end{pmatrix} = \begin{pmatrix} d_x \\ d_y \\ d_z \end{pmatrix} + \mathbf{R} \begin{pmatrix} x \\ y \\ z \end{pmatrix}$$

defined by $3 + 9$ parameters, will not affect the χ^2 of the fits. The matrix can be decomposed into

- three rescaling factors of coordinate axes: f_x, f_y, f_z ,
- three rotations: $\mathcal{D}_x, \mathcal{D}_y, \mathcal{D}_z$ and three shearings: $\mathcal{T}_{xz}, \mathcal{T}_{yz}, \mathcal{T}_{xy}$.

In addition there may be weakly defined **nonlinear** transformations (bend).

These degrees of freedom can be fixed

- by mixture of different data and by external measurements \rightarrow hardware alignment devices, i.e. alignment by tracks has to be supplemented by external information, or fixed
- by equality constraints or by orthogonalization methods.

Equality constraints

Undefined degrees of freedom can be fixed by adding equality constraint equations of the type

$$g(\mathbf{p}) = 0 \quad \text{e.g.} \quad d_x = \sum_i \Delta x_i = 0$$

e.g. “zero average displacement”, or “zero rotation of the whole detector”.

There are several possibilities:

- fix certain parameters (e.g. planes), or
- fix linear combinations of parameters (after diagonalization), or
- add equality constraint equation, i.e. append linearized Lagrange multiplier equation $\lambda (g(\mathbf{p}) + \mathbf{g}^T \cdot \Delta \mathbf{p} = 0)$ with $\mathbf{g} = \partial g(\mathbf{p}) / \partial \mathbf{p}$:

$$\left(\begin{array}{c|c} \mathbf{C}^{\text{global}} & \mathbf{g} \\ \hline \mathbf{g}^T & \mathbf{0} \end{array} \right) \left(\begin{array}{c} \Delta \mathbf{p}^{\text{global}} \\ \lambda \end{array} \right) = \left(\begin{array}{c} \mathbf{b}^{\text{global}} \\ -g(\mathbf{p}) \end{array} \right)$$

Notes: problem with Lagrange function can be solved by matrix inversion, but matrix is no longer positive definit. Alternative is penalty function ($\dots + |g(\mathbf{p})|^2$) or combination of both (augmented Lagrangian).

Parameters of a Vertex detector

Planar sensors (silicon pixel or strip detectors): local (sensor) coordinates $\mathbf{q} = (u, v, w)$ and global detector coordinates $\mathbf{r} = (x, y, z)$ are transformed by

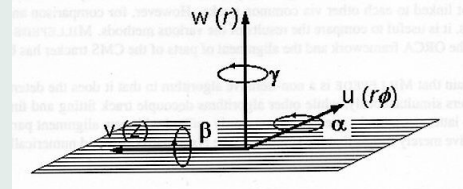
$$\mathbf{q} = \mathbf{R}(\mathbf{r} - \mathbf{r}_0) \quad \mathbf{R} = \text{nominal rotation}, \quad \mathbf{r}_0 = \text{nominal position}$$

After alignment the transformation becomes (with $\Delta\mathbf{q} = (\Delta u, \Delta v, \Delta w)$)

$$\mathbf{q}^{\text{aligned}} = \Delta\mathbf{R} \mathbf{R}(\mathbf{r} - \mathbf{r}_0) - \Delta\mathbf{q}$$

The correction matrix $\Delta\mathbf{R}$ is given by small rotations by α , β and γ around the u -axis, the (new) v -axis and the (new) w -axis:

$$\Delta\mathbf{R} = \mathbf{R}_\alpha \mathbf{R}_\beta \mathbf{R}_\gamma \approx \begin{pmatrix} 1 & \gamma & -\beta \\ -\gamma & 1 & \alpha \\ \beta & -\alpha & 1 \end{pmatrix},$$



where the approximation has sufficient accuracy for small angles α , β and γ .

Six parameters are required for each individual detector element, out of which **three** parameters (two translations, one rotation) are very sensitive.

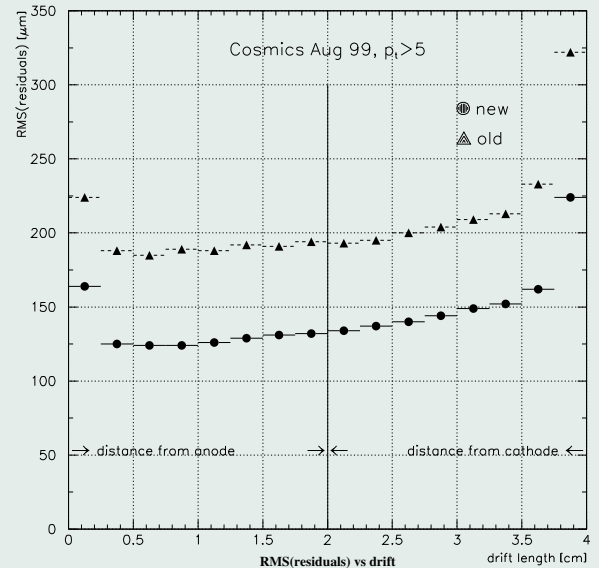
Note: Number of parameters per sensor can be reduced \Rightarrow partially separable functions.

Driftchamber alignment

Claus Kleinwort et al., Detailed calibration of H1 drift chamber using MILLEPEDE with about 1400 parameters.

1400 parameters using 50 000 tracks:

- *common* alignment of the drift chamber and the silicon detector;
- for both CJC1 and CJC2 14 global parameters representing an overall shift or tilt are introduced;
- local variations of the drift velocity v_{drift} for cells halves and layers halves are observed, which are parametrized by $180 + 112$ corrections, which change with the HV configuration;
- for each wire group (8 wires) corrections to T_0 are introduced (330 corrections).



Reduction of residual RMS by MILLEPEDE.

No problems with undefined degrees of freedom for parameters like T_0 or v_{drift} corrections.

First development of Millepede principle (reduction of matrix) 1996 ... and used in H1 1997 ...

V. Blobel: Experience with Online Calibration Methods, Contribution to CHEP'97, Berlin 1997 (including the **Millepede** principle), not accepted.

MILLEPEDE design: experiment-independent program, with well-defined interface to experiment-dependent data.

Available on web page <http://www.desy.de/~blobel>

V. Blobel, Linear Least Squares Fits with a Large Number of Parameters, (2000), 22 pages and full Fortran code.

V. Blobel and C. Kleinwort: A New Method for the High-Precision Alignment of Track Detectors, PHYSTAT2002, Durham, [arXiv-hep-ex/0208021](https://arxiv.org/abs/hep-ex/0208021)

Used (or under test) by H1(1997), CDF(2001), HERA-b, ZEUS, CMS, ATLAS, LHC-b ...??? and rewritten in C++ several times (unpublished).

Principle of reducing matrix size (perhaps) used earlier:

Schreiber, O. (1877): Rechnungsvorschriften für die trigonometrische Abteilung der Landesaufnahme, Ausgleichung und Berechnung der Triangulation zweiter Ordnung. Handwritten notes. Mentioned in W. Jordan (1910): Handbuch der Vermessungskunde, Sechste erw. Auflage, Band I, Paragraph III: 429-433. J.B.Metzler, Stuttgart.

4. Millepede II

Start of development in May 2005 after discussions with Hamburg cms group, with aim:

- alignment with up to 100 000 parameters in a reasonable time on a standard PC;
- keep Millepede principle (simultaneous fit of arbitrary number of tracks and alignment parameters);
- allow different (direct and iterative) methods for the solution of large matrix equation, using mathematical methods from the (mathematical) community (literature) (not home-made iterative methods);
- design with even stronger separation of experiment-dependent program and the Millepede alignment;
- automatic recognition of existing alignment parameters, allowing suppression of parameters with too few data;
- constraints as equality constraints, or like measurements.

Test by PhD student (cms) since summer 2005. Not yet all options realized.

Millepede \Rightarrow Mille + Pede

Mille: small C++ or Fortran routine, called within the experiments event-processing program

Pede: stand-alone experiment-independent alignment program, with many options and input files.

Mille Event loop: write alignment information (derivatives, hits, ...) to special file

Pede Normal equation loop: read data and text files and form C and b

Solution of matrix equation: direct or iterative

solve $C\Delta p = b$ for Δp

End of solution

End of normal equation loop

End of event loop

12 h for event loop, $\frac{1}{2}$ h for ... PEDE II

Event loop: only once to extract the data (more than once, if large non-linearities)

Normal equation loop: once (or repeated, if outlier suppression necessary or for L-BFGS)

Solution of matrix equation: inversion direct ($n < 5000$), or iteratively

Overview over methods

Pede performs (like Millepede I) fits to single e.g. tracks, and generates the vector and matrix of the normal equations in different formats.

No single optimal method, different methods for different conditions (number of parameters, sparsity):

Matrix inversion: ● same routine as in MP I, for up to 5 000 parameters, with time $\propto n^3$;

Diagonalization: ● slower than inversion, allows to recognize insignificant linear combinations (no constraints necessary);

Sparse matrix storage: ● allows to store big sparse matrices

Generalized minimal residual method: ● fast method for large sparse matrices, factor 5 000 faster than inversion for $n = 12\,000$. Routines MINRES ● (and SYMMLQ \ominus);

Preconditioning: \ominus allows to reduce number of iterations, possible in MINRES (and SYMMLQ);

Band matrix decomposition: ★ fast (time $\propto m^2 \cdot n$), but only approximation (but useful for preconditioning?); bordered band matrix possible too;

Limited memory BFGS: ★ uses only *virtual* matrix, low space requirement, but many iterations(?);

Partially separable functions: ★ *trick*, to reduce number of parameters (but more iterations).

Code: ● =included in MILLEPEDE II, ★ = prepared \ominus = not yet tried.

Standard method for symmetric matrices is stable Gauss algorithm with pivot selection of diagonal, with

$$\text{Computing time} = \text{constant} \times n^3 \quad < 1 \text{ hour for } n = \text{several thousand}$$

(in practice “constant” increases for large n).

A standard matrix routine will fail – at least a few parameters out of many thousands will be badly defined.

Subroutine SPMINV (in Millepede) for symmetric matrices in $(n^2 + n)/2$ words: choose largest pivot, but stop inversion if no acceptable pivot found, i.e. invert largest possible submatrix; return zero corrections for remaining parameters.

All variances and covariances available in inverse matrix.

The **global correlation coefficient**, ρ_j is a measure of the total amount of correlation between the j -th parameter and *all* the other variables. It is the largest correlation between the j -th parameter and every possible linear combination of all the other variables.

$$\rho_j = \sqrt{1 - \frac{1}{(\mathbf{V})_{jj} \cdot (\mathbf{C})_{jj}}} \quad \text{and} \quad (\mathbf{V})_{jj} \cdot (\mathbf{C})_{jj} = \frac{1}{1 - \rho_j^2} \quad \mathbf{V} = \mathbf{C}^{-1}$$

Solution by diagonalization

The diagonalization of the symmetric matrix $\mathbf{C} = \mathbf{J}^T \mathbf{J}$ allows to recognize singularity or near singularity by the determination of eigenvalues, and to suppress corresponding linear combinations of parameters.

Algorithms are iterative, computing time ≈ 10 times larger compared to inversion, and solution less precise.

$$\mathbf{C} = \mathbf{U} \mathbf{D} \mathbf{U}^T \quad \text{Diagonalization of symmetric matrix}$$

with \mathbf{D} diagonal, \mathbf{U} square and orthogonal with $\mathbf{U} \mathbf{U}^T = \mathbf{U}^T \mathbf{U} = \mathbf{1}$. Note: $\mathbf{C}^{-1} = \mathbf{U} \mathbf{D}^{-1} \mathbf{U}^T$

eigenvalue ordering in $\mathbf{D} = [\text{diag}(\lambda_i)]$: $\lambda_1 \geq \dots \geq \lambda_k \geq \lambda_{k+1} = \dots \lambda_n = 0$ (or very small)

$$\text{Solution of } \boxed{\mathbf{C} \Delta \mathbf{p} = \mathbf{b}} \quad \text{by} \quad \Delta \mathbf{p} = \mathbf{U} \left[\text{diag} \left(\frac{1}{\sqrt{\lambda_i}} \right) \right] \underbrace{\left[\text{diag} \left(\frac{1}{\sqrt{\lambda_i}} \right) \right] (\mathbf{U}^T \mathbf{b})}_{= \mathbf{q} \quad \text{with} \quad \mathbf{V}[\mathbf{q}] = \mathbf{1}}$$

with $1/\lambda_i = 0$ for $\lambda_i = 0$ or small q_i with $|q_i| \lesssim 1$

\Rightarrow Suppression of insignificant linear combinations, which could produce distortions of the detector.

Singular value decomposition avoids the formation of normal equations ($\mathbf{C} = \mathbf{J}^T \mathbf{J}$) and is numerically more accurate than normal-equation methods.

$$\mathbf{J} = \mathbf{V} \mathbf{D} \mathbf{U}^T \quad \text{Singular value decomposition (SVD) for } m \times n \text{ matrix } \mathbf{J}$$

with \mathbf{D} diagonal, \mathbf{U} square and orthogonal with $\mathbf{U} \mathbf{U}^T = \mathbf{U}^T \mathbf{U} = \mathbf{1}$ and $m \times n$ matrix \mathbf{V} column-orthogonal with $\mathbf{V}^T \mathbf{V} = \mathbf{1}$. Diagonal elements σ_i of \mathbf{D} are called singular values, with $\sigma_i^2 = \lambda_i$.

$$\text{Solution of } \boxed{\min \|\mathbf{J} \Delta \mathbf{p} - \mathbf{r}\|_2} \quad \text{by} \quad \Delta \mathbf{p} = \mathbf{U} \left[\text{diag} \left(\frac{1}{\sigma_i} \right) \right] (\mathbf{V}^T \mathbf{r})$$

with $1/\sigma_i = 0$ for $\sigma_i = 0$ or close to 0.

⇒ Suppression of insignificant linear combinations, which could produce distortions of the detector.

Advantage: numerically accurate due to small condition number $\kappa(\mathbf{D}) = \sqrt{\kappa(\mathbf{J}^T \mathbf{J})}$

Disadvantage: requires to store huge $m \times n$ matrices \mathbf{J}/\mathbf{V} and large $n \times n$ matrix \mathbf{U} , and can **not** be used (in this form) in large alignment problems. SVD can be applied to the matrix $\mathbf{C} = \mathbf{J}^T \mathbf{J}$; then it is equivalent to diagonalization.

Sparse matrix storage

Large matrices are usually *sparse*, with small fraction of non-zero off-diagonal elements (fraction called q). Mathematical methods for the solution of matrix equations exist, which only require the product of the matrix with vectors.

The indexed storage scheme of PCGPACK, modified for a symmetric matrix, is used: it requires arrays with

$$n + q \cdot n(n - 1)/2 \quad \text{double precision (data) and integer (indices) words}$$

and is optimized for the product (9 lines of code).

During matrix generation (sums): a (binary) search is necessary to find the location for an index pair (i, j)

The automatic generation of [parameter-index](#) relations requires a large number of comparisons.

Example: For 100 000 parameters and 1 Mio tracks, each 20 hits about $2q \times 10^{19}$ sequential comparisons would be necessary, for $q = 1\%$ this corresponds to **6 cpu-years** (1 nanosec/comparison).

A faster method using a combination of hashing, sorting and binary search is used.

Generalized minimal residual method (GMRES)

Solution of a very large system of linear equations with sparse matrix, by an clever solution of a quadratic minimization problem, in analogy to the method of conjugate gradients (Hestenes, Stiefel 1952).

Example: MINRES (M. A. Saunders), designed to solve

system of linear equations

$$\boxed{\mathbf{C} \Delta \mathbf{p} = \mathbf{b} \quad \text{or} \quad \min \|\mathbf{C} \Delta \mathbf{p} - \mathbf{b}\|_2}$$

where \mathbf{C} is a symmetric matrix of logical size $n \times n$, which may be indefinite and/or singular, very large and sparse. It is accessed *only* by means of a subroutine call

call Aprod (n, x, y) to return $\mathbf{y} = \mathbf{C}\mathbf{x}$

for any given vector \mathbf{x} .

Example of **compact storage**: row-index sparse storage *)

For $n = 100\,000$ and 512 Mbytes memory: $q = 1\%$

5 Gbytes memory: $q = 10\%$

C. C. Paige and M. A. Saunders (1975), Solution of sparse indefinite systems of linear equations, SIAM J. Numer. Anal. 12(4), pp. 617-629.

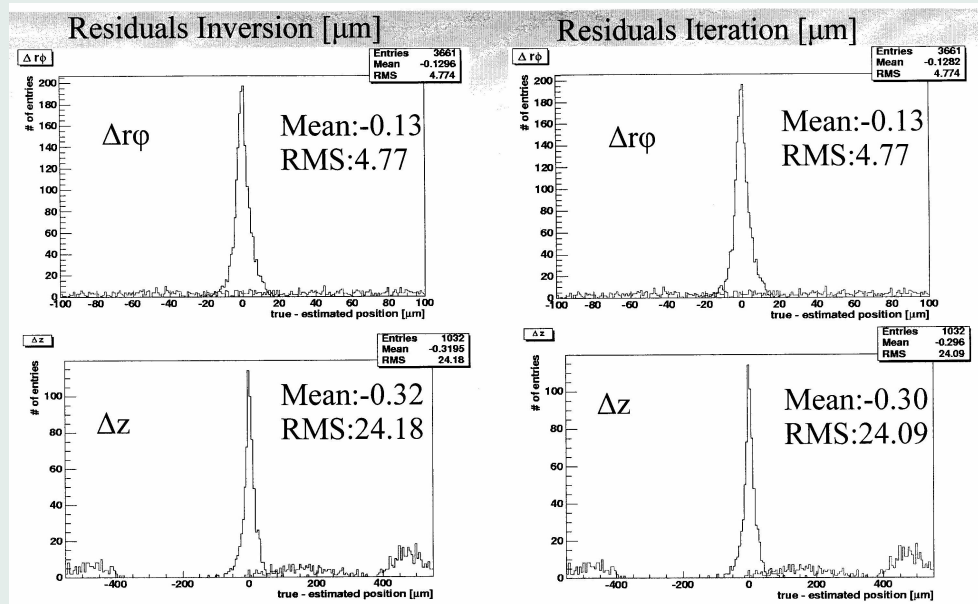
www.stanford.edu/group/SOL/software/minres.html

*) W. H. Press, S. A. Teukolsky, W. T. Vetterling, B.P. Flannery, NUMERICAL RECIPES – The Art of Scientific Computing, Cambridge Univ.

Comparison

For 12 000 parameters

- matrix inversion (cpu-time 12 h, 46 min, 5 s), and
- iterative solution with MINRES (cpu-time 32 s).



Example: 250 000 tracks, 500 Mbytes file, 4 400 variable parameters, 3.1 % non-zero off-diagonal elements (plot \Rightarrow ps-file), 100 sec for preparation, and 100 sec per iteration (250 000 track fits + solution).

Preconditioning

The convergence rate of iterative methods depends on the spectral properties of the matrix.

A transformation may improve the spectral properties:

$$\begin{array}{ll} \text{instead of} & \mathbf{C} \Delta \mathbf{p} = \mathbf{b} \\ \text{solve} & \mathbf{M}^{-1} \mathbf{C} \Delta \mathbf{p} = \mathbf{M}^{-1} \mathbf{b} \quad \text{for } \Delta \mathbf{p}, \end{array}$$

which has the same solution as the original system, but the condition number of $\mathbf{M}^{-1} \mathbf{C}$ may be smaller.

The matrix \mathbf{M}^{-1} should be an approximation to \mathbf{C}^{-1} , such that $\mathbf{M}^{-1} \mathbf{C} \approx \mathbf{1}$.

One possibility: band portion of the inverse of the band matrix approximation of matrix \mathbf{C} .

Preconditioning is an option in MINRES by means of a subroutine call

$$\text{call Msolve(n, x, y) ,} \quad \text{to solve } \mathbf{M} \mathbf{y} = \mathbf{x} \quad \text{for } \mathbf{y}$$

without altering vector \mathbf{x} .

Cholesky decomposition

Decomposition of the symmetric matrix \mathbf{C}

$$\mathbf{C} = \mathbf{L}\mathbf{D}\mathbf{L}^T \quad \text{decomposition}$$

is “numerically extremely stable”, and can be made *in-space*. Matrix \mathbf{L} is a left unit triangular matrix (diagonal elements =1) and \mathbf{D} is a diagonal matrix (less stable for semi-definite matrix).

Solution of $\mathbf{C}\Delta\mathbf{p} = \mathbf{L}(\mathbf{D}\mathbf{L}^T\Delta\mathbf{p}) = \mathbf{b}$ by

$$\mathbf{L}\mathbf{v} = \mathbf{b} \quad \text{forward substitution} \quad \mathbf{L}^T\Delta\mathbf{p} = \mathbf{D}^{-1}\mathbf{v} \quad \text{forward substitution}$$

With clever ordering (?) of parameters the matrix \mathbf{C} can be approximated by a band matrix.

Band matrices with band-width m : the band structure is kept in this decomposition and computing time for solution is only $\propto n \times m^2$. Selected elements of the inverse matrix (corresponding to the band of the original matrix \mathbf{C}) can be calculated quickly.

Fast methods exist also for

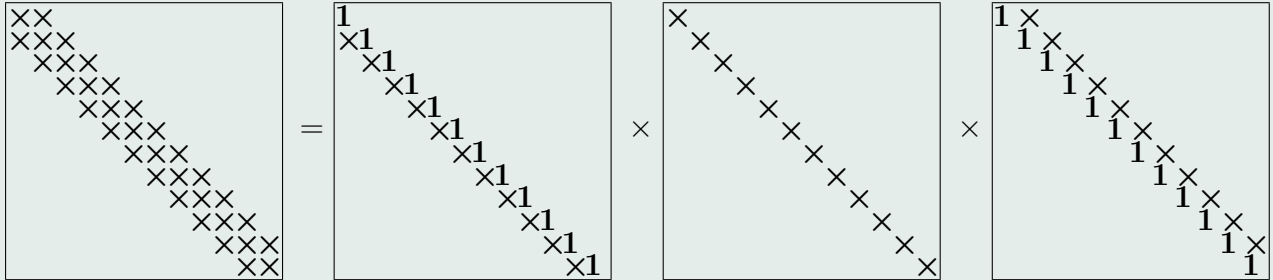
Variable-bandwidth matrices (‘sky-line’ matrix), and for

Bordered band matrices (‘arrow’ matrix), where e.g. border (additional full rows and columns) is due to Lagrange multiplier constraints.

Decomposition is similar to the original Cholesky factorization; the introduction of the diagonal matrix \mathbf{D} avoids the square-root operations.

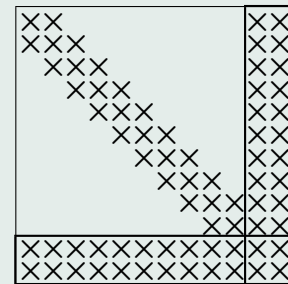
Band matrix decomposition

$$C = L D L^T$$



Note: the inverse of a bandmatrix is a full dense matrix, but
the bandpart of the inverse can be calculated fast.

In reality there will be *bordered* band matrices: \Rightarrow
even for those cases there are fast *direct* methods for
the solution – and for the determination of elements
of the inverse matrix.



Quasi-Newton methods ...

...require only the gradient of the objective function $F(\mathbf{p})$ at each iteration. New step $\Delta\mathbf{p}$ for line search of $F(\mathbf{p} + \alpha \cdot \Delta\mathbf{p})$ is calculated from

$$\mathbf{H} \Delta\mathbf{p} = \nabla F \quad \text{or} \quad \mathbf{B} \nabla F = \Delta\mathbf{p} \quad \text{with } \mathbf{B} = \mathbf{H}^{-1}$$

with approximate Hessian \mathbf{H} or inverse Hessian \mathbf{B} .

Starting from a simple assumption (e.g. $\gamma\mathbf{1}$) the Hessian \mathbf{H} or the inverse Hessian \mathbf{B} is improved by updates, using the difference vectors

$$\mathbf{s}_k = \mathbf{p}_{k+1} - \mathbf{p}_k \quad \mathbf{y}_k = \nabla F_{k+1} - \nabla F_k .$$

Requiring the secant equation $\mathbf{H}_{k+1}\mathbf{s}_k = \mathbf{y}_k$ or $\mathbf{B}_{k+1}\mathbf{y}_k = \mathbf{s}_k$, the most-popular update formula is

$$\rho_k = 1/\mathbf{y}_k^T \mathbf{s}_k \quad \mathbf{B}_{k+1} = (\mathbf{1} - \rho_k \mathbf{s}_k \mathbf{y}_k^T) \mathbf{B}_k (\mathbf{1} - \rho_k \mathbf{y}_k \mathbf{s}_k^T) + \rho_k \mathbf{s}_k \mathbf{s}_k^T \quad (\text{BFGS})$$

To minimize $F(\mathbf{p} + \alpha \cdot \Delta\mathbf{p})$ the line-search has to satisfy the Wolfe (or strong Wolfe) conditions.

The BFGS method has $\mathcal{O}(n^2)$ operations per iteration and has a superlinear rate of convergence ... but requires of course to store the full (dense) matrix \mathbf{B} .

Quasi-Newton method, a revolutionary idea, invented by physicist W.C. Davidon (Argonne) in the mid 1950s; paper not accepted for publication. Different update formulas were developed during the following 20 years.

Limited memory BFGS

One step of the BFGS method has the form

$$\mathbf{p}_{k+1} = \mathbf{p}_k - \alpha \cdot \mathbf{B}_k \nabla F_k \quad (\text{line search, } \alpha \text{ usually close to } 1)$$

The limited memory BFGS (short: L-BFGS) method uses update information only [from most recent iterations.](#)

With $\mathbf{B}_0 = \gamma \mathbf{1}$ the product $\mathbf{B}_k \nabla F_k$ can be evaluated from the last m difference vectors $\mathbf{y}_k, \mathbf{s}_k$, a “memoryless” BFGS method.

- *Good* values for m are in the range $3 \dots 20$;
- the storage requirement with $n(2m + 4)$ is linear in n , and
- $\approx 3/2 m^2 n$ operations are needed per iteration.
- The rate of convergence is *linear*.

... an optimization method for $n \gg 100\,000$ parameters ?

J. Nocedal, *Updating quasi-Newton matrices with limited storage*, Mathematics of Computation **35** (1980), 773 - 782.

D.C. Liu and J. Nocedal, *On the limited-memory BFGS method for large scale optimization*, Mathematical Programming **45** (1989) 503 - 528

Partially separable functions

In *separable* optimization, the objective function $F(\mathbf{p})$ is decomposed into a sum of simpler functions that can be optimized separately, for example

$$F(\mathbf{p}) = F_1(\mathbf{p}_1) + F_2(\mathbf{p}_2) + F_3(\mathbf{p}_3)$$

...requires detailed information about the structure of the objective function.

Can the method be applied to alignment?

Consider the six parameters $\Delta\mathbf{p}$ of a single sensor (translation + rotation) and the corresponding 6×6 matrix \mathbf{C} and right-hand side 6-vector \mathbf{b} :

$$\mathbf{C} \Delta\mathbf{p} = \mathbf{b} \quad \text{with diagonalization} \quad \mathbf{C} = \mathbf{U} \mathbf{D} \mathbf{U}^T$$

The transformed vector of linear combinations

$$\mathbf{q} = \mathbf{U}^T \mathbf{b}$$

has a diagonal covariance matrix $\mathbf{V} = \mathbf{D}^{-1}$, and thus the linear combinations are uncorrelated.

Idea: consider only the three most-sensitive (i.e. with largest diagonal elements) or most-significant linear combinations, or try two pairs of each three linear combinations, which are *almost* independent. This reduces the storage space of the matrix to 1/4.

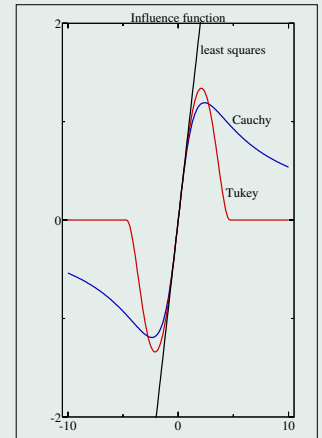
The presence of outliers in the data can deteriorate the alignment result. Difficulty: wrong initial alignment parameters can fake outliers.

Millepede I: Large initial cut at $\approx 10\sigma$ reduced to 3σ in ≈ 5 iterations.

Millepede II: No cut in first iteration, followed by technique of M-estimates in subsequent iterations.

The objective function in least squares is the sum of **squares** of scaled residuals z , with **larger influence for larger residuals** (outliers). The **square** is replaced in M-estimates by a dependence with reduced influence for larger residuals.

	$\rho(z) = \ln \text{pdf}(z)$	influence function $\psi(z) = d\rho(z)/dz$	add. weight $\omega(z) = \psi(z)/z$
Least squares	$= \frac{1}{2} z^2$	$= z$	$= 1$
Cauchy($c = 2.3849$)	$= \frac{c^2}{2} \ln \left(1 + (z/c)^2 \right)$	$= \frac{z}{1 + (z/c)^2}$	$= \frac{1}{1 + (z/c)^2}$
Huber $\begin{cases} \text{if } z \leq c = 1.345 \\ \text{if } z > c = 1.345 \end{cases}$	$= \begin{cases} z^2/2 \\ c(z - c/2) \end{cases}$	$= \begin{cases} z \\ c \cdot \text{sign}(z) \end{cases}$	$= \begin{cases} 1 \\ c/ z \end{cases}$



Summary

A comparison of large scale optimization methods

Residuals	$F(\mathbf{p})$	∇F	$\mathbf{H} \equiv \nabla^2 F$	\mathbf{H}^{-1}	Method	MP
×	—	—	—	—	Residual	
—	×	×	—	—	Steepest descent	
—	×	×	—	×	Quasi-Newton (e.g. BFGS)	
—	×	×	×	×	Newton (standard optimization)	
—	—	×	×	×	Inversion	I,II
—	—	×	×	×	Diagonalization	II
—	—	×	(×)	—	Sparse \mathbf{H} solution	II
—	—	×	(×)	(×)	Band or bordered band matrix	II
—	×	×	—	(×)	L-BFGS (virtual \mathbf{H}^{-1})	II

- Several solution modes in MILLEPEDE II should allow alignment even with $n = 100\,000$ parameters.
- Options to add constraints and to reduce influence of outliers included.
- Simultaneous use of several (all) types of events and data – physics data: single tracks, vertices, invariant-mass constraints, and cosmes, overlaps . . . external measurements is essential for accurate alignment.

Contents

1. Introduction	2	Parameters of a Vertex detector	27
Optimization	3	Driftchamber alignment	28
Time	4	History	29
Space	5	4. Millepede II	30
Iterations and convergence	6	Decay of Millepede	31
Rate of convergence	7	Overview over methods	32
Least squares	8	Solution by matrix inversion	33
2. Alignment of a toy detector	9	Solution by diagonalization	34
First attempt based on residuals	10	Solution by singular value decomposition	35
Result from the first attempt	11	Sparse matrix storage	36
First attempt – Discussion	12	Generalized minimal residual method (GMRES)	37
Results from the second attempt	13	Comparison	38
Use of higher mathematics?	14	Preconditioning	39
Results from a simultaneous fit	15	Cholesky decomposition	40
Determination of drift velocities	16	Band matrix decomposition	41
A more realistic scenario	17	Quasi-Newton methods	42
3. Millepede I	18	Limited memory BFGS	43
Normal equations of least squares	19	Partially separable functions	44
Normal equations	20	M-estimates	45
Reduction of matrix size	21	Summary	46
Numerical linear algebra	22	Table of contents	47
Solution by matrix inversion	23		
Global correlations	24		
Undefined degrees of freedom	25		
Equality constraints	26		