

Alignment Algorithms

Volker Blobel – Universität Hamburg

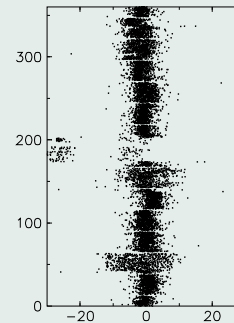
Abstract

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. Many different methods are used in HEP, which either require a large number of iterations, due to a linear convergence behaviour, or which require the solution of a large matrix equation. 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.

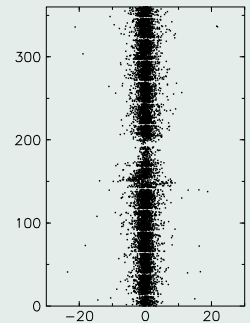
1. Introduction and classification
2. Alignment of a toy detector
3. Millepede I and II
4. Mathematical methods

Summary

Before alignment ...



... after alignment



1. Introduction and classification

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, improved mass and vertex resolution, 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 .

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 on survey data, perhaps on Laser alignment data, and **mainly** on track residual minimization (*incomplete* data, several degrees of freedom undefined);
- no “reference standards with known values” exist (\approx exceptions are data from $e^+e^- \rightarrow \mu^+\mu^-$ and cosmics with $\mathbf{B} = 0$).

Alignment parameter corrections Δp

Geometrical corrections: six parameters are required for each individual detector element or larger detector component: translations: $(\Delta u, \Delta v, \Delta w)$ and angles: (α, β, γ) .

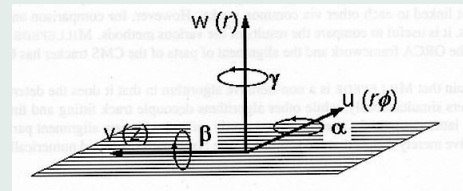
Often three parameters (two translations, one rotation) are very sensitive.

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

$$(u, v, w) = \mathbf{R}(\mathbf{r} - \mathbf{r}_0) \quad \mathbf{R}, \mathbf{r}_0 = \text{nominal rotation, position}$$
$$(u, v, w)_{\text{aligned}} = \Delta \mathbf{R} \mathbf{R}(\mathbf{r} - \mathbf{r}_0) + (\Delta u, \Delta v, \Delta w)$$

The correction matrix $\Delta \mathbf{R}$ is given by small rotations by α, β, 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.

Drift chamber: Lorentz-angle, T_0 , drift velocity v_{drift} , ... global, and per plane, per layer ... and e.g. correction dependent on angle between track and drift direction.

Optimization

Alignment as an optimization problem:

- definition of an objective function $F(\Delta\mathbf{p})$, to be minimized;
- large number n of correlated parameters, up to $n \approx 6 \times 5\,000$ (Atlas) or $n \approx 6 \times 20\,000$ (CMS);
- non-linearities are expected to be not large.

The standard minimization method to obtain $\Delta\mathbf{p}$ = corrections to alignment parameters requires the solution of a system of linear equations:

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

(\mathbf{C} is a n -by- n matrix) or of a sequence of such systems of linear equations in case of non-linearities.

Most (not all) physicists think:

- *The solution of a matrix equation requires a matrix inversion!* **not true!**
- *Matrix inversion is impossible for a large matrix!* **not true!**

“The solution of 4200 equations in 4200 unknowns is computationally infeasible. Even worse, non-linear fit won’t converge.”

Therefore most alignment methods work only with small matrices.

Classification

- **Minimization of sum of squares of residuals:** A “global” objective function $F(\dots)$ or χ^2 -function is constructed

$$F(\dots) = \sum_{\text{data sets}} \left(\sum_{\text{events}} \left(\sum_{\text{tracks}} \left(\sum_{\text{hits}} \Delta_i^2 / \sigma_i^2 \right) \right) \right)$$

Biased algorithms: methods with solution of many small system of equations; requires iterations to reduce or remove bias.

$F(\Delta\mathbf{p}) = \dots$ depends on the alignment parameters corrections $\Delta\mathbf{p}$.

Almost all alignment methods tried in HEP experiments are of this type.

Unbiased algorithms: direct methods with solution of large system of equations, allowing to apply constraints.

$F(\Delta\mathbf{p}, \mathbf{q}) = \dots$ depends on the corrections to alignment parameters $\Delta\mathbf{p}$ and **all** track parameters \mathbf{q}

Examples: Millepede (used by several experiments), Brueckman de Renstroem (Atlas)

- **Special methods** (discussed in other talks)

Kalman filter: \Rightarrow Talk by Rudi Frühwirth: “Alignment using a Kalman Filter Technique”

Update of alignment parameters (sequential parameter estimation) during track processing.

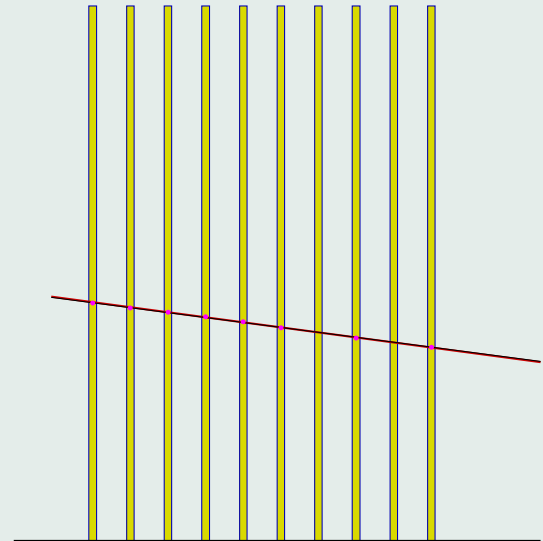
SLD: \Rightarrow Talk by Fred Wickens: “Alignment experience from SLD”

Fits of different functional forms to various residual types, and extraction of alignment parameters from fitted functional forms.

2. Alignment of a toy detector

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

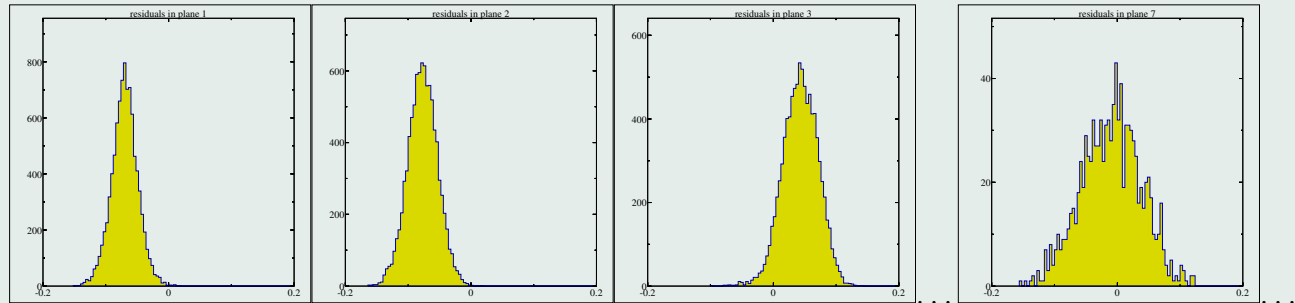
- 10 planes of tracking chambers, 1 m high, 10 cm distance, no magnetic field;
- accuracy $\sigma \approx 200\mu\text{m}$, with efficiency $\epsilon = 90\%$;
- plane 7 sick: accuracy $\sigma \approx 400\mu\text{m}$, with efficiency $\epsilon = 10\%$;
- 10 000 tracks with 82 000 hits available for alignment;
- **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, ignoring the misalignment.
- The residuals (fitted coordinate f_i – measured vertical) are histogrammed, separately for each plane.



- The mean value of the residuals is taken as correction to the vertical plane position.
- One iteration reduces the residuals significantly, ... but fitted track parameters are essentially unchanged (biased).

This is the standard method used iteratively 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.

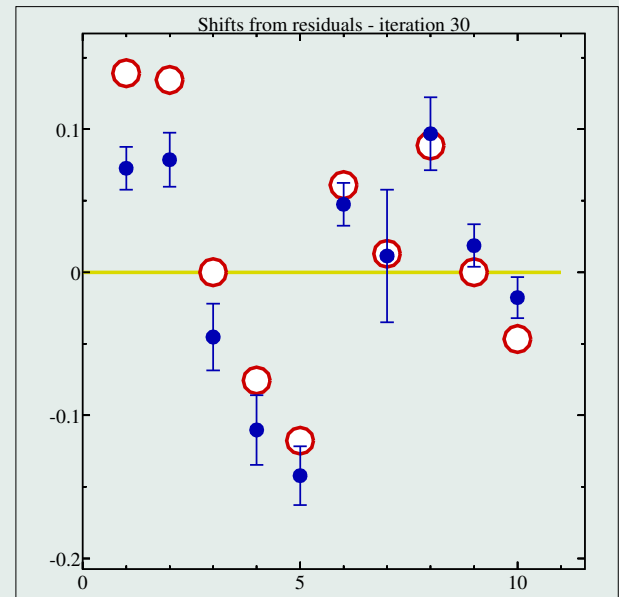
After 30 iterations ...

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

Units are cm.

red circle = true shift (displacement)

blue disc = displacement, determined from residuum



First attempt – Discussion

The result is not (yet) encouraging!

Method is equivalent to attempt, to solve a linear system of 10 unknowns, which is singular.

The reason for non-convergence is simple:

Two degrees of freedom are undefined: a simultaneous shift and an 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.

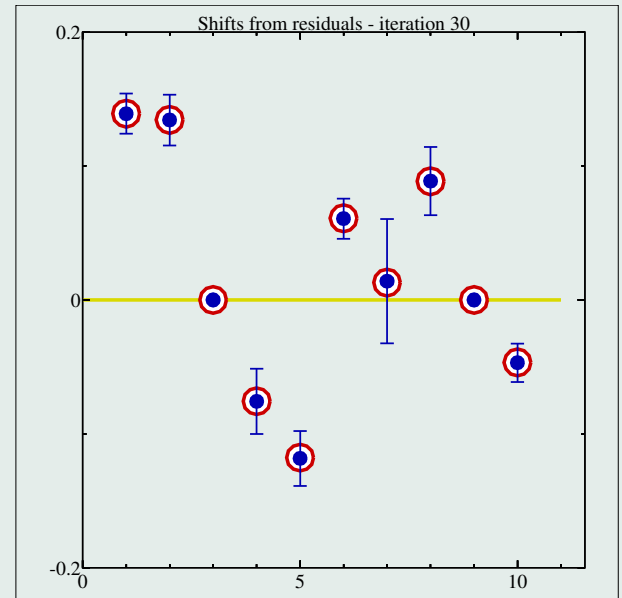
After 30 iterations with planes 3 and 9 fixed (displacement = 0) ...

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

Units are cm.

red circle = true shift (displacement)

blue disc = displacement, determined from residuum



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

$$\text{linear convergence} \quad \frac{\|\mathbf{x}_{k+1} - \mathbf{x}^*\|}{\|\mathbf{x}_k - \mathbf{x}^*\|} \leq r \quad r \in (0, 1)$$

for all k sufficiently large. Speed of convergence depends on the eigenvalue spectrum – slow convergence for small eigenvalues.

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:

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

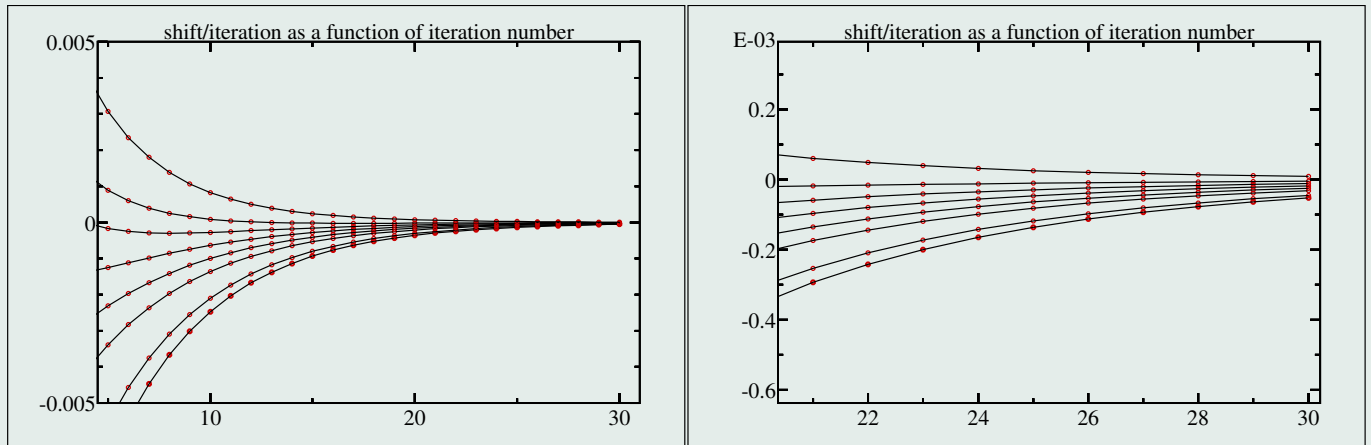
Often in applications small $\|\mathbf{x}_{k+1} - \mathbf{x}_k\|$ is assumed to indicate convergence, but

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

i.e. a small $\|\mathbf{x}_{k+1} - \mathbf{x}_k\|$ does not mean small distance to the solution!

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

Many small shifts $\mathbf{x}_k - \mathbf{x}^*$ add up during the iterations.



... linear convergence.

Units are cm.

Alternative: unbiased algorithm

Include the alignment parameters in the parameters fitted in track fits – requires a simultaneous fits of many tracks, with simultaneous 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 a 20010-by-20010 matrix.

The MILLEPEDE principle allows to reduce the least squares problem to $n = 10$ parameters, which has a direct unbiased solution (no iterations)!

Results from a simultaneous fit

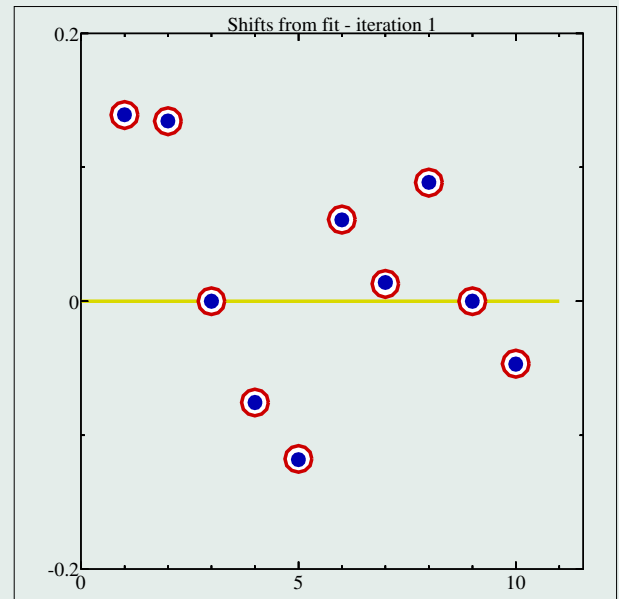
After one step (with planes 3 and 9 fixed at displacement = 0) ...

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, units are cm.)

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: $\Delta v_{\text{drift}}/v_{\text{drift}}$

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

$a_j^{\text{global}} = \text{shift for plane } j$

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

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.

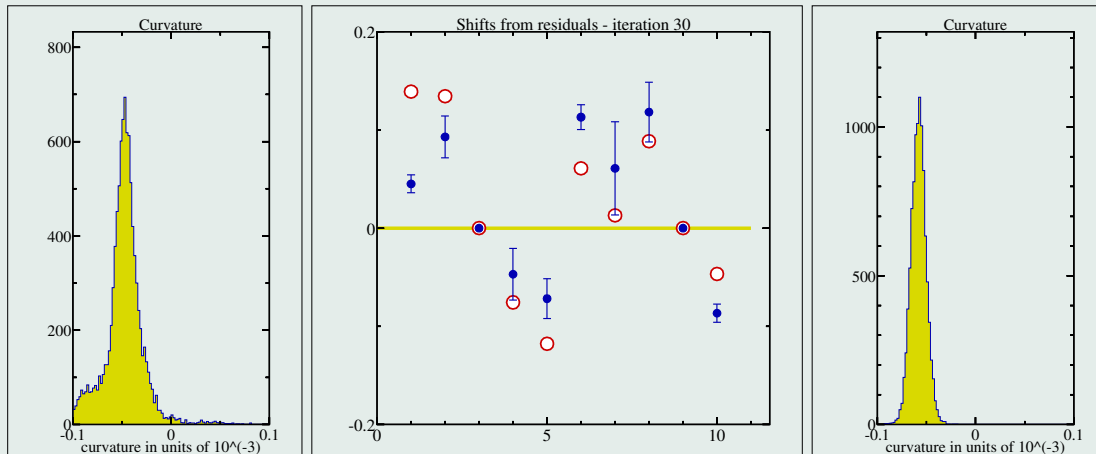
The next improvement would be the introduction of wire T_0 's – additional $10 \times 25 = 250$ parameters.

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:

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

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

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.

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

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

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 $\mathbf{H}_k^{\text{global-local}}$

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

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

Solution by partitioning

Partitioning the matrix equation $\mathbf{C}\mathbf{a} = \mathbf{b}$ (\mathbf{C} symmetric, \mathbf{C}_{11} and \mathbf{C}_{22} are square matrices):

$$\left(\begin{array}{c|c} \mathbf{C}_{11} & \mathbf{C}_{12} \\ \hline \mathbf{C}_{12}^T & \mathbf{C}_{22} \end{array} \right) \begin{pmatrix} \mathbf{a}_1 \\ \mathbf{a}_2 \end{pmatrix} = \begin{pmatrix} \mathbf{b}_1 \\ \mathbf{b}_2 \end{pmatrix}$$

If $\mathbf{C}_{12} \equiv 0$: $\mathbf{C}_{22} \mathbf{a}_2^* = \mathbf{b}_2$ $\mathbf{a}_2^* = \mathbf{C}_{22}^{-1} \mathbf{b}_2$ (called *local* solution)

Submatrix \mathbf{B} of the complete inverse matrix corresponding to the upper left part \mathbf{C}_{11} is easy to calculate afterwards: $\mathbf{B} = (\mathbf{C}_{11} - \mathbf{C}_{12}\mathbf{C}_{22}^{-1}\mathbf{C}_{12}^T)^{-1}$ and complete inverse matrix equation in terms of \mathbf{B} :

$$\begin{pmatrix} \mathbf{a}_1 \\ \mathbf{a}_2 \end{pmatrix} = \begin{pmatrix} \mathbf{B} & -\mathbf{B}\mathbf{C}_{12}\mathbf{C}_{22}^{-1} \\ \hline -\mathbf{C}_{22}^{-1}\mathbf{C}_{12}^T\mathbf{B} & \mathbf{C}_{22}^{-1} - \mathbf{C}_{22}^{-1}\mathbf{C}_{12}^T\mathbf{B}\mathbf{C}_{12}\mathbf{C}_{22}^{-1} \end{pmatrix} \begin{pmatrix} \mathbf{b}_1 \\ \mathbf{b}_2 \end{pmatrix}$$

Solution for subvector \mathbf{a}_1 : $(\mathbf{C}_{11} - \mathbf{C}_{12}\mathbf{C}_{22}^{-1}\mathbf{C}_{12}^T) \mathbf{a}_1 = (\mathbf{b}_1 - \mathbf{C}_{12}\mathbf{a}_2^*)$

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 covariance matrix $\mathbf{V}_k = (\mathbf{C}_k^{\text{local}})^{-1}$ of the local parameters
- 2. Derivatives:** calculate for all hits (index i) the vectors of derivatives $\boldsymbol{\delta}_i^{\text{local}}$ and $\mathbf{d}_i^{\text{global}}$ for all local and the relevant global parameters, and update matrices:

$$\mathbf{C} := \mathbf{C} + \sum_i w_i \mathbf{d}_i^{\text{global}} (\mathbf{d}_i^{\text{global}})^T \quad \mathbf{b} := \mathbf{b} + \sum_i w_i r_i \mathbf{d}_i^{\text{global}} \quad \mathbf{H}_k = \sum_i w_i \mathbf{d}_i^{\text{global}} (\boldsymbol{\delta}_i^{\text{local}})^T$$

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 two '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:

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

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

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

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

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

Subroutine SPMINV (in Millepede); in-space inversion of symmetric matrix in $(n^2 + n)/2$ words: choose always 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.

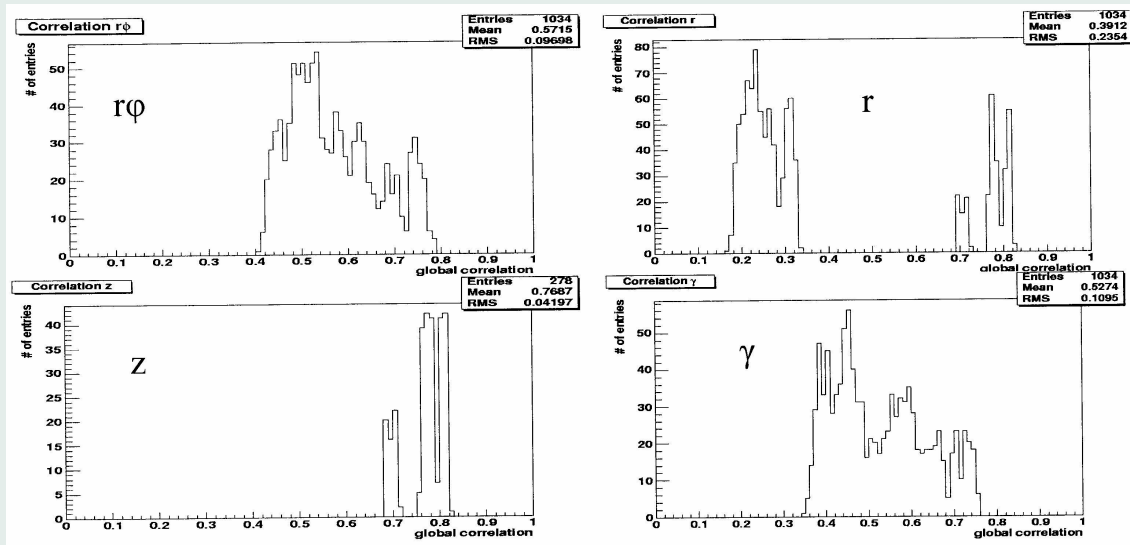
The **global correlation coefficient**, ρ_j is a measure of the total amount of correlation between the j -th parameter and *all* the other variables. It is the largest correlation between the j -th parameter and every possible linear combination of all the other variables.

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

Matrix is ill-conditioned (almost singular), if any ρ close to 1, with large condition number κ .

Global correlations

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



Values ≈ 1 means strong correlation and almost singular matrix – inversion may be impossible (and biased iterative methods would be extremely slow).

Values depend on geometry and type of data – additional data (cosmics, vertex and mass-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).

Degrees of freedom can be fixed, if necessary,

- 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 e.g. by fixed planes;
- by equality constraints (e.g. $d_x = 0$) 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 $\mathbf{g} = \partial g(\mathbf{p}) / \partial \mathbf{p}$ (determination of stationary point of Lagrange function, matrix not positive definite.):

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

- alternative is penalty function ($\dots + |g(\mathbf{p})|^2$) or combination of both (augmented Lagrangian);
- elimination method, with reduction of the total number of parameters.

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 I, code 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)

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

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

⇒ Talk by Claus Kleinwort: combined alignment/calibration of H1 vertex detector and drift chamber (resolution improved by factor 2).

Formalism equiv. to Millepede principle derived for ATLAS:

P.Bruckman, S.Haywood, Least Squares Approach to the Alignment of the Generic High Precision Tracking System, PHYSTAT05, Oxford 12-15 Sept. 2005

Principle of reducing matrix size (perhaps) used already in 19.th century in surveying.

Millepede II

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

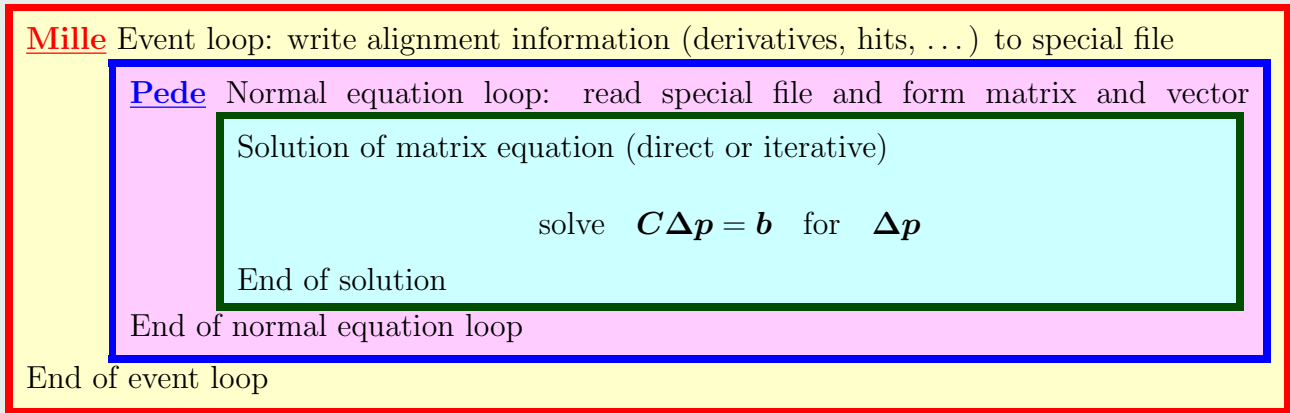
- alignment with up to 100 000 parameters in a reasonable time on a standard PC;
- keep Millepede principle (unbiased, simultaneous fit of arbitrary number of tracks and of alignment parameters);
- allow different (direct and iterative) methods for the solution of large matrix equation, using mathematical methods from the (mathematical) community (literature) (no home-made iterative methods);
- design with even stronger separation of experiment-dependent program and the Millepede alignment computation;
- 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



12 h for event loop, $\frac{1}{2}$ h for MILLEPEDE 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 ($n < 5000$), or iterative solution (MINRES, or L-BFGS)

The introduction of overall equality constraints requires the solution of large systems of equations!

How to solve very large systems of equations?

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

Matrix inversion: ● e.g. routine 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); possible for large n on special hardware;

Sparse matrix storage: ● allows to store big sparse matrices

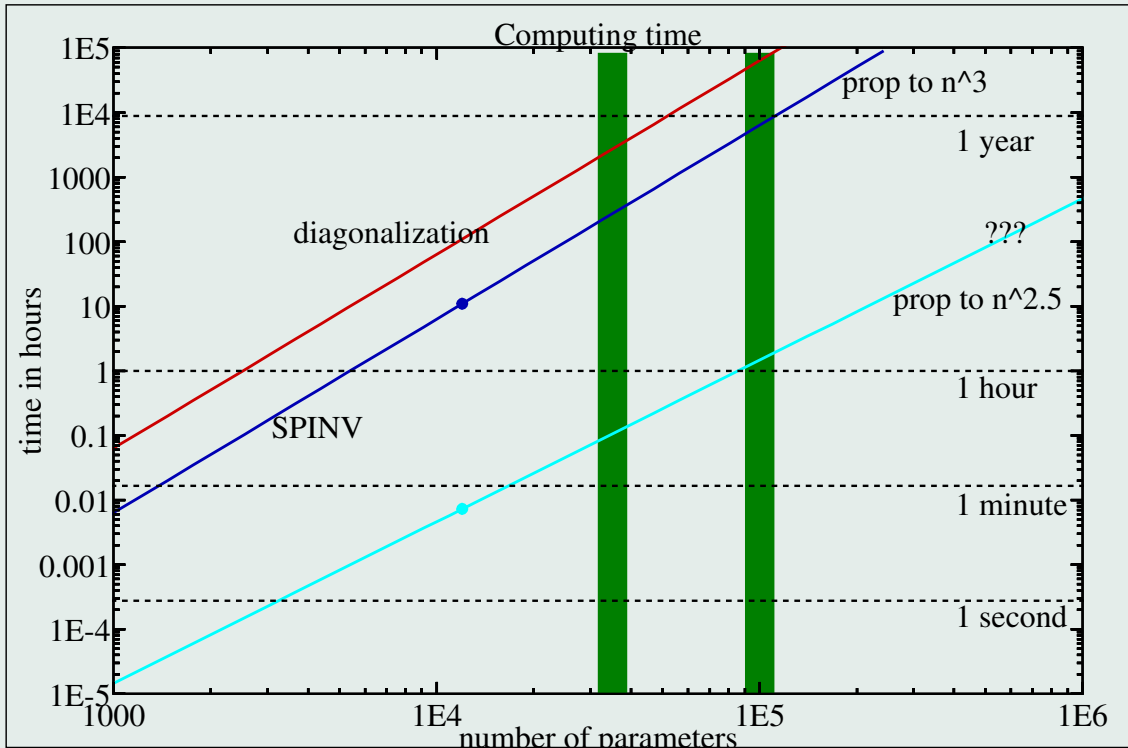
Generalized minimal residual method: ● fast method for large sparse matrices, factor $> 1\,000$ faster than inversion for $n = 12\,000$. Routines MINRES ● (and SYMMLQ ⊖);

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

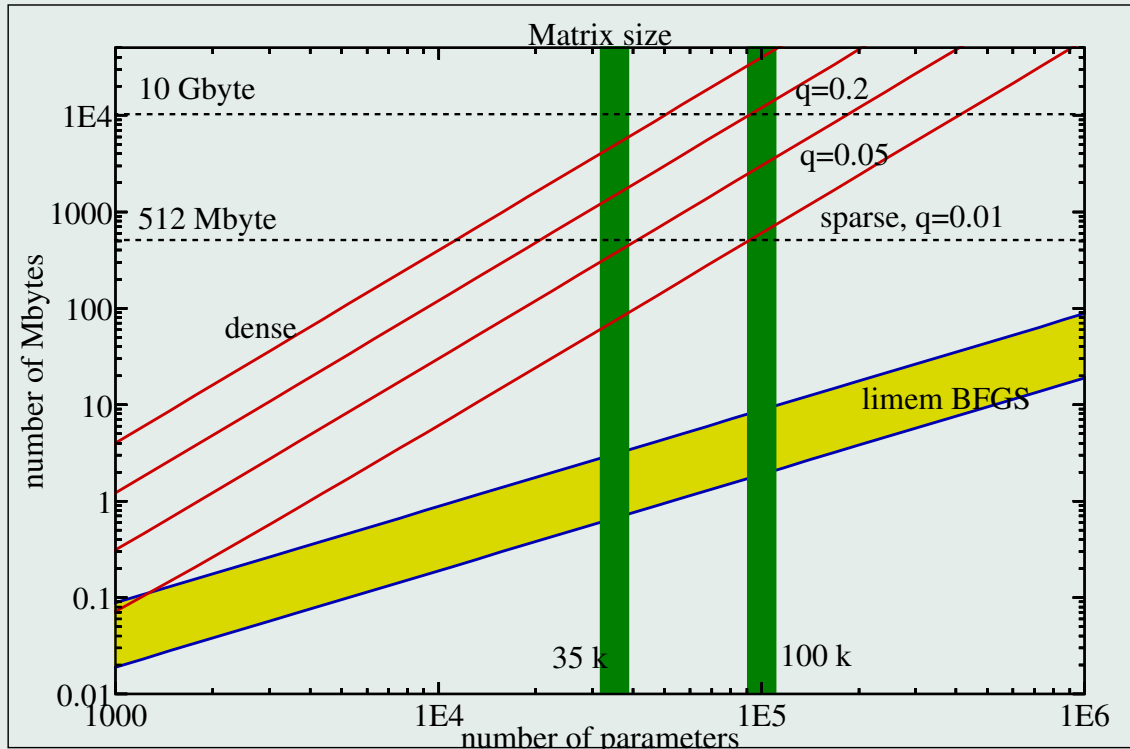
Limited memory BFGS: ● uses only *virtual* matrix, low space requirement, but many iterations(?);

Millepede II Code: ● =included, ⊖ = not yet tried.

Method of M-estimates instead of cuts against outliers; square (of least squares) replaced by density with larger tails for outliers.



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



q = fraction of non-zero off-diagonal elements

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 replacement $1/\lambda_i = 0$ 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.

Reduction of the number of parameters

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

$$\mathbf{C}_{6 \times 6} \Delta\mathbf{p} = \mathbf{b} \quad \text{with diagonalization} \quad \mathbf{C}_{6 \times 6} = \mathbf{U}\mathbf{D}\mathbf{U}^T$$

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.

(2) One (small) set of alignment parameters for large detector section, added to single-sensor alignment parameters which are either *active* (with $\sum = 0$ constraint), or *non-active* (fixed parameters).

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 $\sigma_i^2 = \lambda_i$.

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

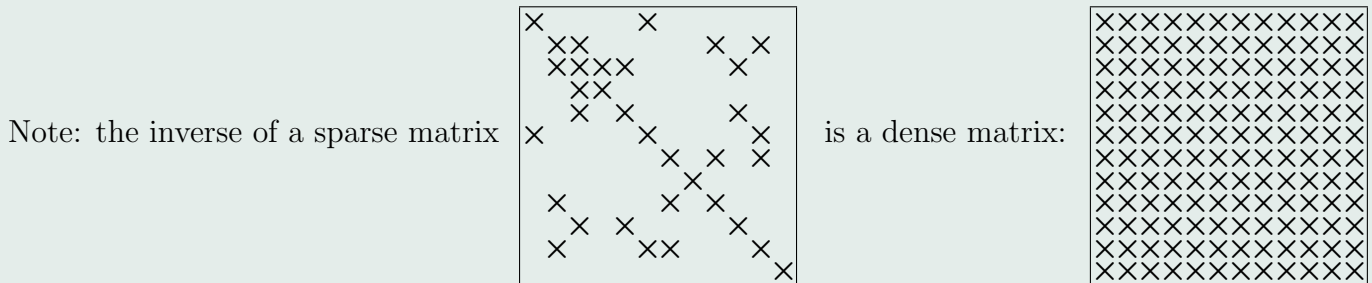
with $1/\sigma_i = 0$ for $\sigma_i = 0$ or close to 0. \Rightarrow **Suppression of insignificant linear combinations, which could produce distortions of the detector.**

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 q of non-zero off-diagonal elements. The automatic generation of [parameter-index](#) relations requires a large number of comparisons. A fast method using a combination of hashing, sorting and binary search is used in MP II.



Mathematical methods for the solution exist, which only require the product of the matrix with vectors: $\mathbf{y} = \mathbf{C}\mathbf{x}$.

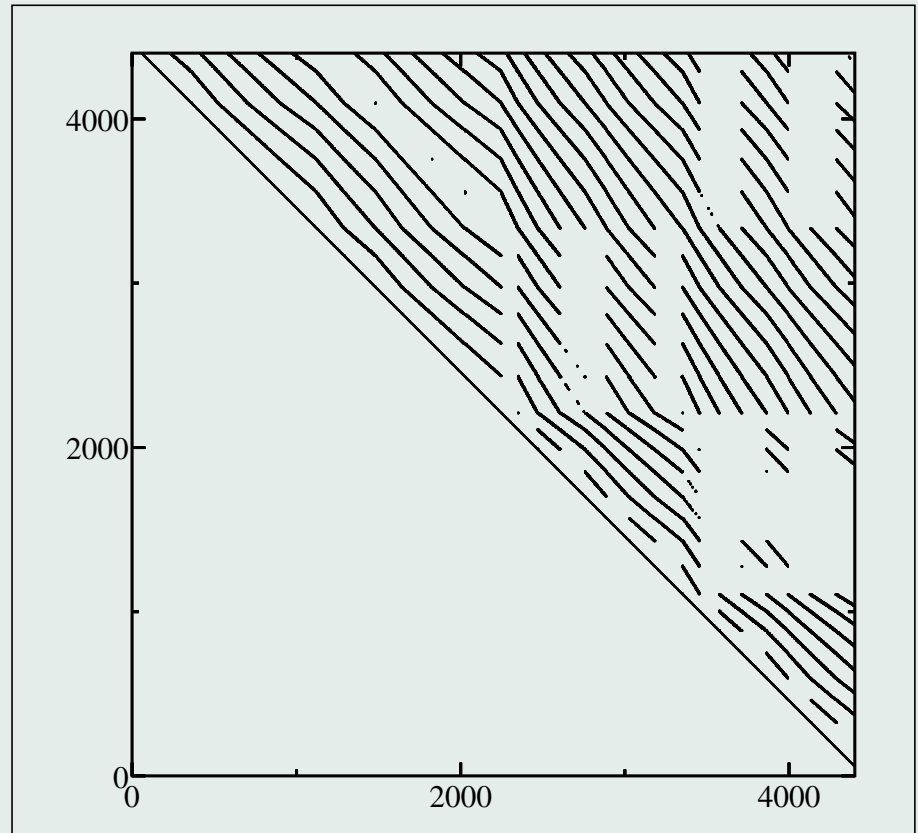
The indexed storage scheme of PCGPACK, modified for a *symm.* matrix, is used: it requires arrays with

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

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)

A sparse matrix example

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



Generalized minimal residual method (GMRES)

Solution of a very large system of linear equations with sparse matrix, by an optimized 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, 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} .

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

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

without altering vector \mathbf{x} . The matrix \mathbf{M}^{-1} should be an approximation to \mathbf{C}^{-1} , such that $\mathbf{M}^{-1}\mathbf{C} \approx \mathbf{1}$. (not yet implemented).

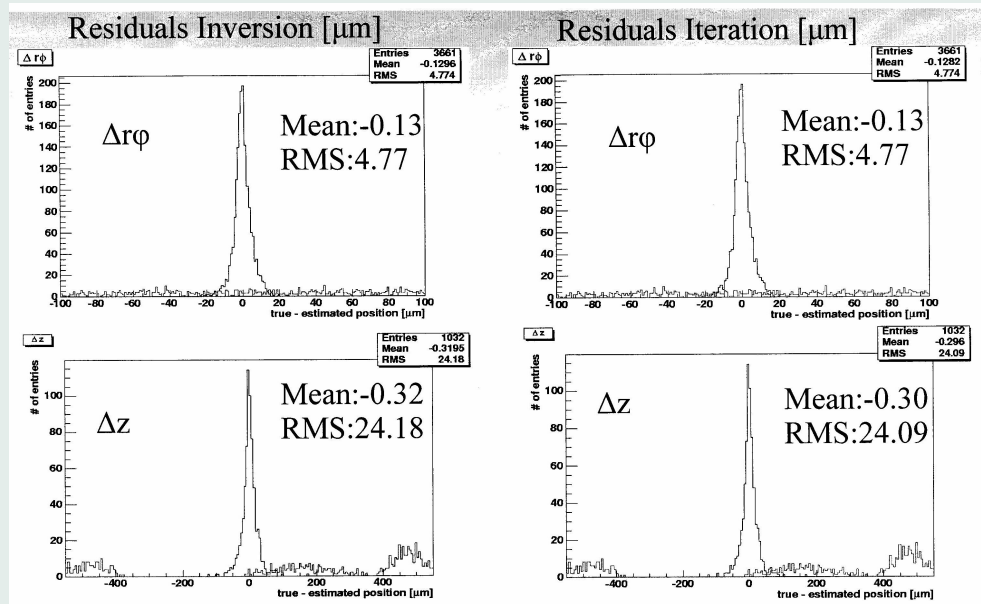
C. C. Paige and M. A. Saunders (1975), Solution of sparse indefinite systems of linear equations, SIAM J. Numer. Anal. 12(4), pp. 617-629.

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

Comparison

For 12 000 parameters

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



Quasi-Newton methods . . .

. . . require only the gradient ∇F 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

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

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

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

(used e.g. in MINUIT/MIGRAD).

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} and thus cannot be used.

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 the m 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 $\mathbf{y}_k, \mathbf{s}_k$ only (no matrix storage).

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

... 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

Jorge Nocedal and Stephen J.Wright, Numerical Optimization, Springer

Summary

Two classes of alignment algorithms are developed in HEP experiments:

Biased algorithm: require a certain number of iterations and requires to fix certain planes (in 3D);

Unbiased algorithm: requires to solve a large system of equations, allows to add equality constraints in formalism.

Both methods reduce average residuals – what about behaviour for long-range correlations between alignment parameters ?

Solution methods, in order of increasing size parameter n :

$F(\mathbf{p})$	∇F	$\mathbf{H} \equiv \nabla^2 F$	$\mathbf{V} \equiv \mathbf{H}^{-1}$	Method
–	×	×	×	Diagonalization
–	×	×	×	Inversion
–	×	(×)	–	Sparse \mathbf{H} solution
×	×	–	(×)	L-BFGS (virtual \mathbf{H}^{-1})

All these solution methods are implemented in MILLEPEDE II, which should become able to perform alignment even with $n \approx 100\,000$ parameters.

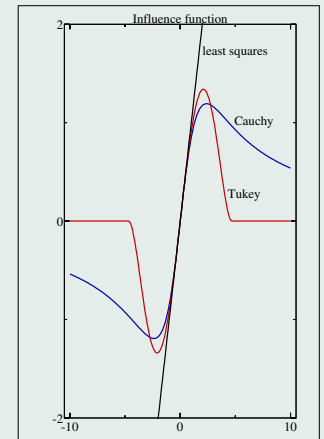
The presence of outliers in the data can deteriorate the alignment result. Difficulty: wrong initial alignment parameters can fake outliers.

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 $\begin{cases} \text{if } z \leq c = 1.345 \\ \text{if } z > c = 1.345 \end{cases}$	$= \begin{cases} z^2/2 \\ c(z - c/2) \end{cases}$	$= \begin{cases} z \\ c \cdot \text{sign}(z) \end{cases}$	$= \begin{cases} 1 \\ c/ z \end{cases}$



Covariance matrix with MINRES

The inverse of matrix \mathbf{C} is the covariance matrix \mathbf{V} of the alignment parameters. This is available with matrix inversion and diagonalization, but not with MINRES.

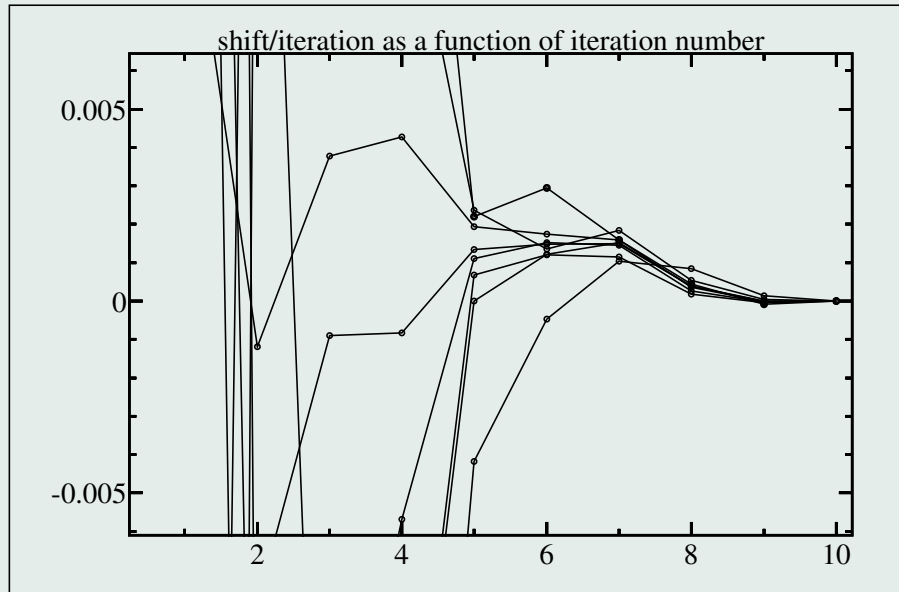
Method to compute some elements of \mathbf{V} with MINRES:

Solution of matrix equation $\mathbf{C}\mathbf{V} = \mathbf{1}$ right hand-side is unit matrix
for \mathbf{V} would give the complete covariance matrix \mathbf{V} .

Solution of matrix equation $\mathbf{C}\mathbf{v} = \mathbf{u}$ right hand-side is column-vector of unit matrix
for \mathbf{V} will give on column of the covariance matrix \mathbf{V} .

Elements of covariance matrix are determined by hits statistics and by geometry.

Fast convergence for limited memory BFGS method after iteration 8.



... fast rate of linear convergence.

Units are cm.

Contents

	“Decay” of Millepede ...	27	
1. Introduction and classification	2	4. Mathematical methods	28
Alignment parameter corrections Δp	3	Solution time	29
Optimization	4	Matrix space	30
Classification	5	Solution by diagonalization	31
2. Alignment of a toy detector	6	Reduction of the number of parameters	32
First attempt based on residuals	7	Solution by singular value decomposition	33
Result from the first attempt	8	Sparse matrix storage	34
First attempt – Discussion	9	A sparse matrix example	35
Results from the second attempt	10	Generalized minimal residual method (GMRES)	36
Iterations and convergence	11	Comparison	37
Shift per iteration	12	Quasi-Newton methods	38
Alternative: unbiased algorithm	13	Limited memory BFGS	39
Results from a simultaneous fit	14	Summary	40
Determination of drift velocities	15	M-estimates	41
A more realistic scenario	16	Covariance matrix with MINRES	42
3. Millepede	17	Shift per iteration	43
Normal equations	18	Table of contents	44
Solution by partitioning	19		
Reduction of matrix size	20		
Solution by matrix inversion in Millepede I	21		
Global correlations	22		
Undefined degrees of freedom	23		
Equality constraints	24		
History	25		
Millepede II	26		