

# Bundle Adjustment Math Reading Seminar Thing

LukášP

Brno University of Technology  
Faculty of Information Technology  
Božetěchova 2, 612 66 Brno  
[ipolok@fit.vutbr.cz](mailto:ipolok@fit.vutbr.cz)



- 3D geometry
- Projective geometry
- Problem Statement
- Minimum Geometric Problems
- Least squares, Moore-Penrose pseudoinverse, SVD
- Nonlinear Least Squares
- Intermezzo I – computing derivatives
  - Symbolic / Numerical differentiation
  - Dual numbers
  - Lie groups basics
- Handling Severe Nonlinearity
  - Levenberg Marquardt
  - Dogleg
- Handling outliers – robust estimation
- Intermezzo II – solving linear systems
  - Direct methods
  - Schur complement trickery
  - ~~• Iterative methods~~
  - ~~• Inexact step NLS~~
- Calculating covariances

- Position representations
  - Euclidean  $[x, y, z]$  3D
  - Inverse depth  $[x/z, y/z, 1/z]$  3D
  - Inverse distance  $[x, y, z, 1/d]$ , where  $\|[x, y, z]\| = 1$  1D
- Rotation representations
  - Rotation matrix  $[r, u, f]$  9D
    - Hard to constrain orthogonality in numerical manipulation
  - (Unit) Quaternion  $[x, y, z] \cdot \sin(\theta/2), \cos(\theta/2)$  4D
    - Double cover of  $SO(3)$ !
  - Axis-angle  $[x, y, z] \cdot \theta$  3D
  - Exponential map of  $\mathfrak{so}(3)$   $[x, y, z] \cdot \theta$  3D

- Use homogenous coordinates
  - $[x/w, y/w, z/w] \rightarrow [x, y, z, w]$
- Projection matrix

• Vision

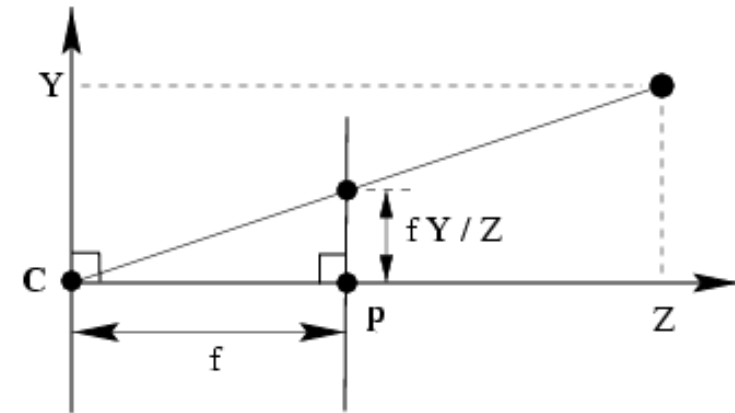
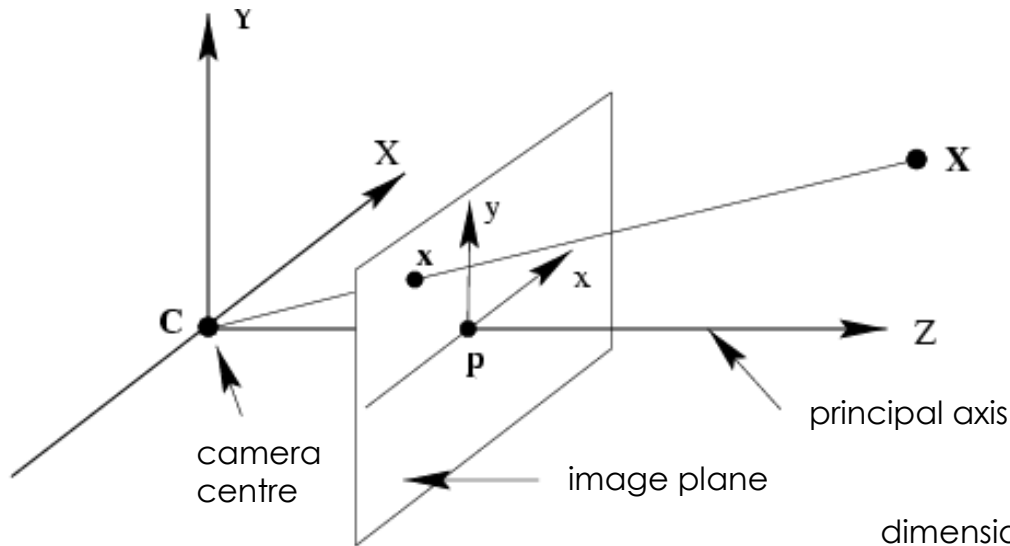
$$\begin{pmatrix} f & s & c_x \\ & f & c_y \\ & & 1 \\ & & 1 & 0 \end{pmatrix}$$

We're not interested in the  $w$  coordinate (it is „always“ 1).  
Leads to division by depth.

• OpenGL

$$\begin{pmatrix} \frac{2n}{r-l} & & \frac{r+l}{r-l} \\ & \frac{2n}{t-b} & \frac{t+b}{t-b} \\ & & \frac{n+f}{n-f} & \frac{2nf}{n-f} \\ & & & -1 \end{pmatrix}$$

- Stolen from Bronek Pribyl's VGE lecture (I think)



intrinsics K

dimension wrangling magic

pose

$$\lambda \begin{bmatrix} x \\ y \\ 1 \end{bmatrix} = \begin{bmatrix} f & p_x & 1 \\ 0 & f & p_y \\ 0 & 0 & 1 \end{bmatrix} \begin{bmatrix} X \\ Y \\ Z \\ 1 \end{bmatrix}$$

$\mathbf{P} = \mathbf{K} [\mathbf{R} \ \mathbf{t}]$

$\lambda \mathbf{x} = \mathbf{P} \cdot \mathbf{X}$

- Spherical
- Curvilinear
- Catadioptric
- ...

- Brown's radial distortion model

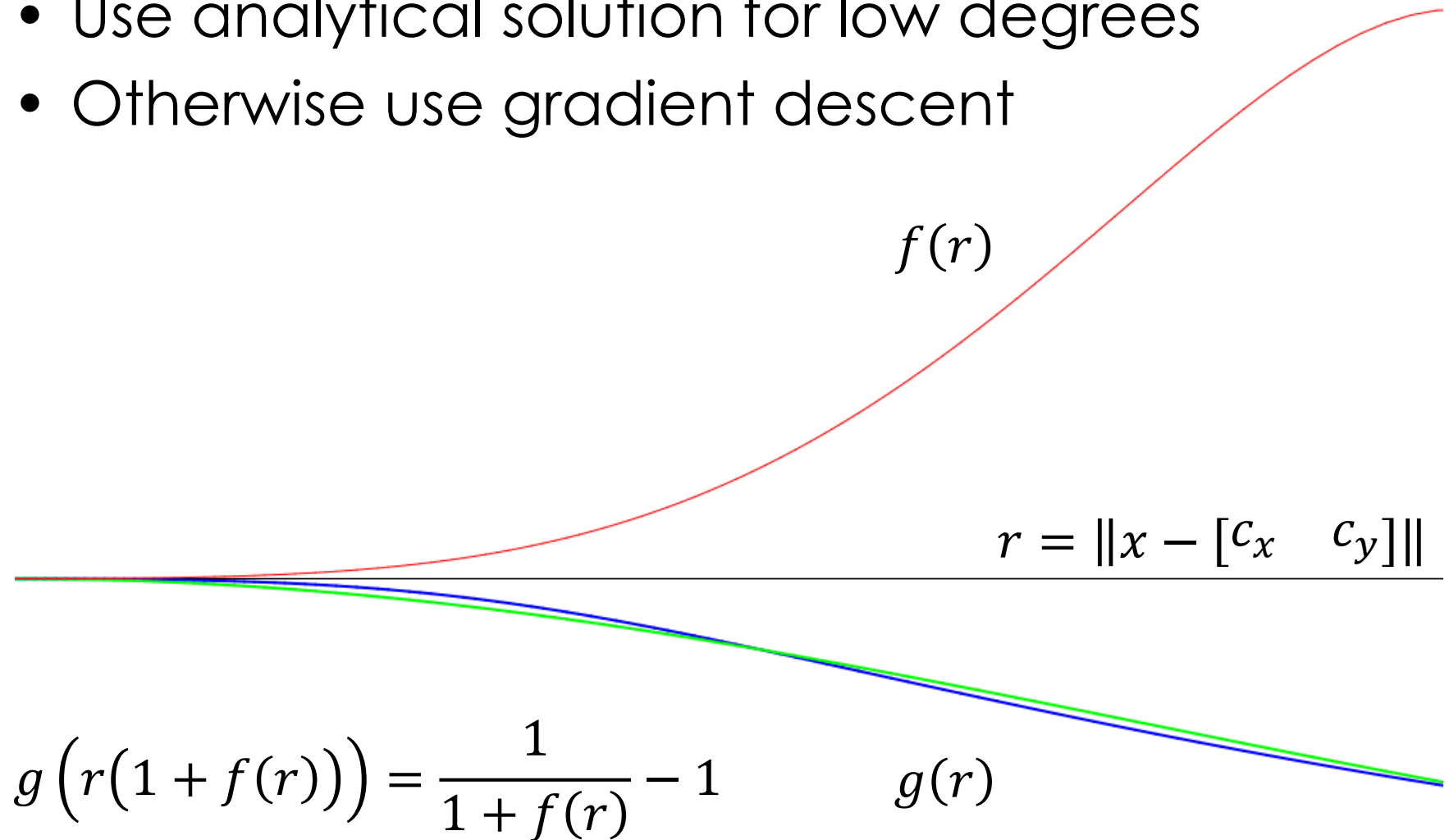
$$x_d = x(1 + f(\|x - [c_x \ c_y]\|))$$

- Here,  $f(\cdot)$  is typically a polynomial function, e.g.:

$$f(a) = [a^3 \ a^5 \ a^7]^T \mathbf{f} ,$$

where  $\mathbf{f}$  is a vector of polynomial coefficients

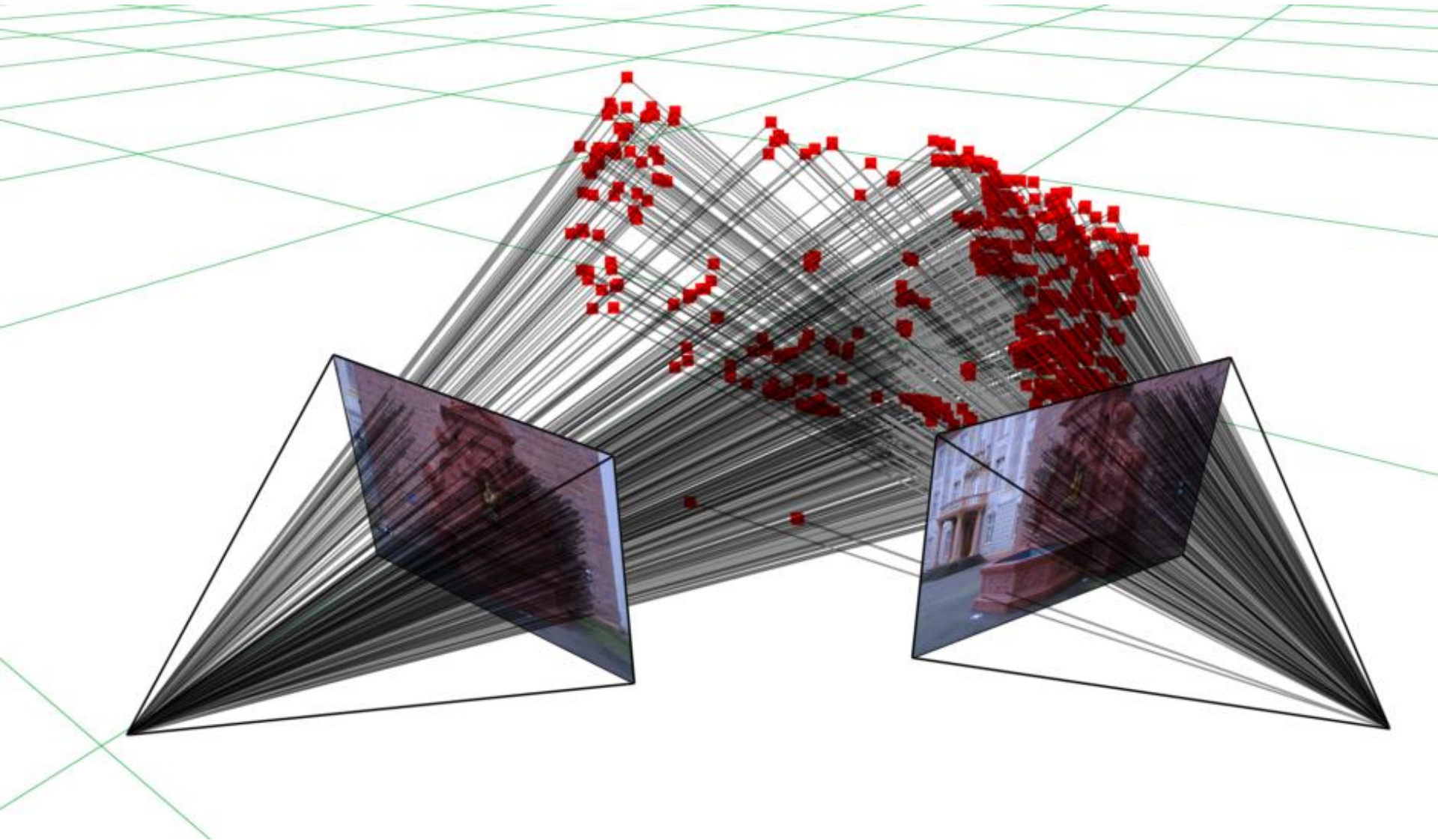
- Problem with invertibility of polynomial  $f(\cdot)$
- Use analytical solution for low degrees
- Otherwise use gradient descent





- Bundle Adjustment
  - „Having a lot of images of a scene, let's build a 3D model“
  - Sparse approach (bundler)
    - Find interest point in the images
    - Match the interest points
    - Calculate spatial rotation / translation between images
    - Triangulate points
  - Semi-dense approach (LSD)
    - Somehow triangulate points in the first camera pair
    - Color the points (from image pixels)
    - For following cameras, reproject the 3D points, optimize rotation / translation to minimize color difference
  - Dense approach (PMVS / CMVS)

- Bundle Adjustment
  - „Having a lot of images of a scene, let's build a 3D model“
  - Sparse approach (bundler)
  - Semi-dense approach (LSD)
  - Dense approach (PMVS / CMVS)
    - Solve dense pixel correspondences
    - Early methods modifications of dynamic programming
    - Gives per-pixel dense depth
      - Requires known relative camera poses (ok for stereo)
    - Can solve alignment by e.g. Iterative closest point (ICP)



- How do we calculate positions from images?
- MGPs! <http://cmp.felk.cvut.cz/mini/>  
(also recent VGS-IT by Tom Pajdla)
- Given a bare minimum of data, get solution by applying some geometric constraints

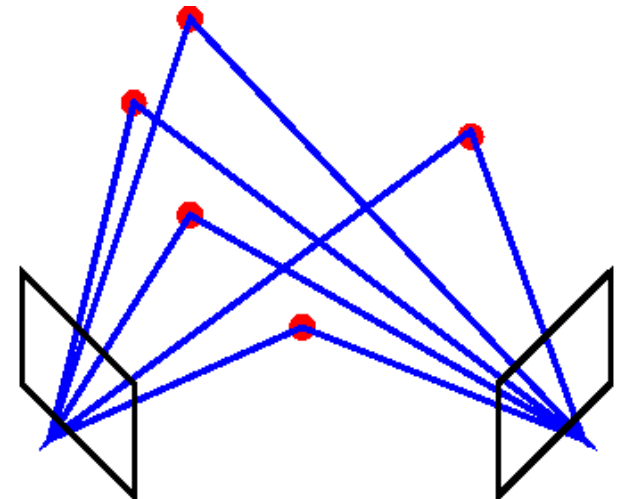
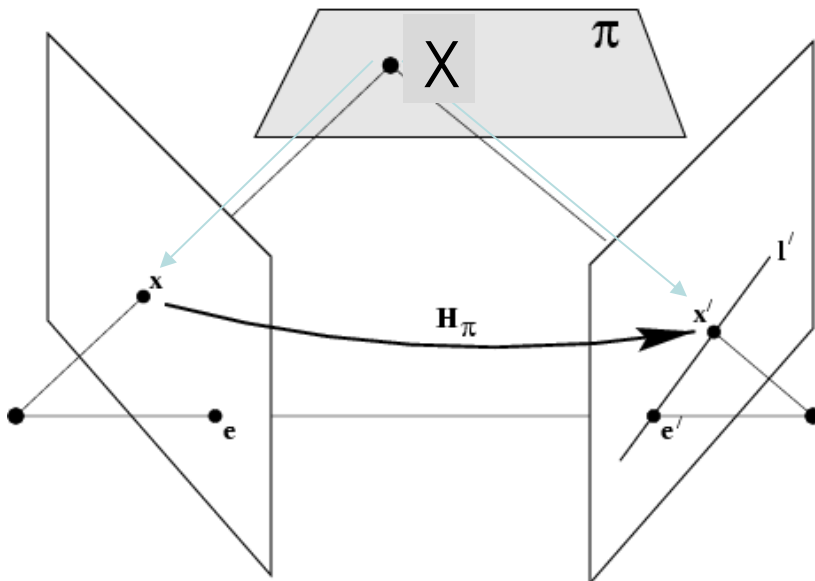
- Five-point algorithm [Nister et al. 2004]
  - Given five 2D points, find  $[R \ t]$  between the cameras
  - Employs Epipolar geometry
    - Fundamental matrix
    - Epipolar constraint
    - Solve 10<sup>th</sup> order polynomial
    - Decompose  $E$  to  $R$  and  $t$

$$F = K_2^{-T}([t]_{\times}R)K_1^{-1} = K_2^{-T}EK_1^{-1}$$

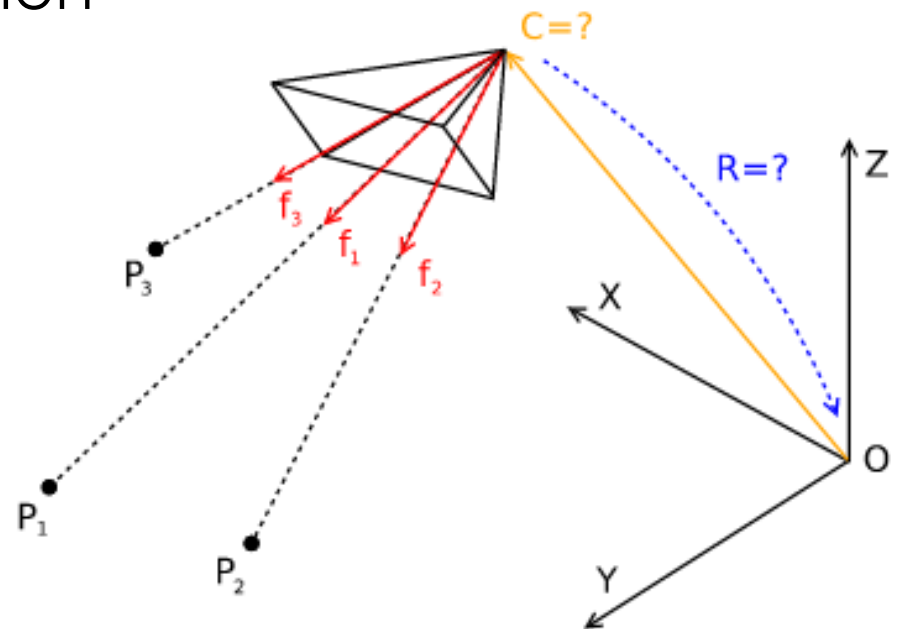
$$x'Fx = 0$$

meh

multiple (im)possible solutions



- Perspective three Point (P3P) [Kneip et al. 2013]
  - Given three 3D points and their 2D observations, find  $[R \ t]$  of the observing camera
  - From knowledge of  $K$ , convert 2D observations to 3D directions
  - Leads to quartic equation
  - Up to four solutions



- Kabsch's algorithm [Kabsch 1976]
  - Many variants (Coustias, Horn, Umeyama, ...)
  - Given two sets of 3D points, calculate  $[R \ t]$  that aligns them
  - Calculate centroids  $c_1 = \frac{1}{N} \sum_{i=1}^N x_i$  and  $c_2 = \frac{1}{N} \sum_{i=1}^N x'_i$
  - Translation  $t = c_2 - c_1$
  - Rotation derived from covariance

$$A = \sum_{i=1}^N (x_i - c_1)^T (x'_i - c_2)$$

as

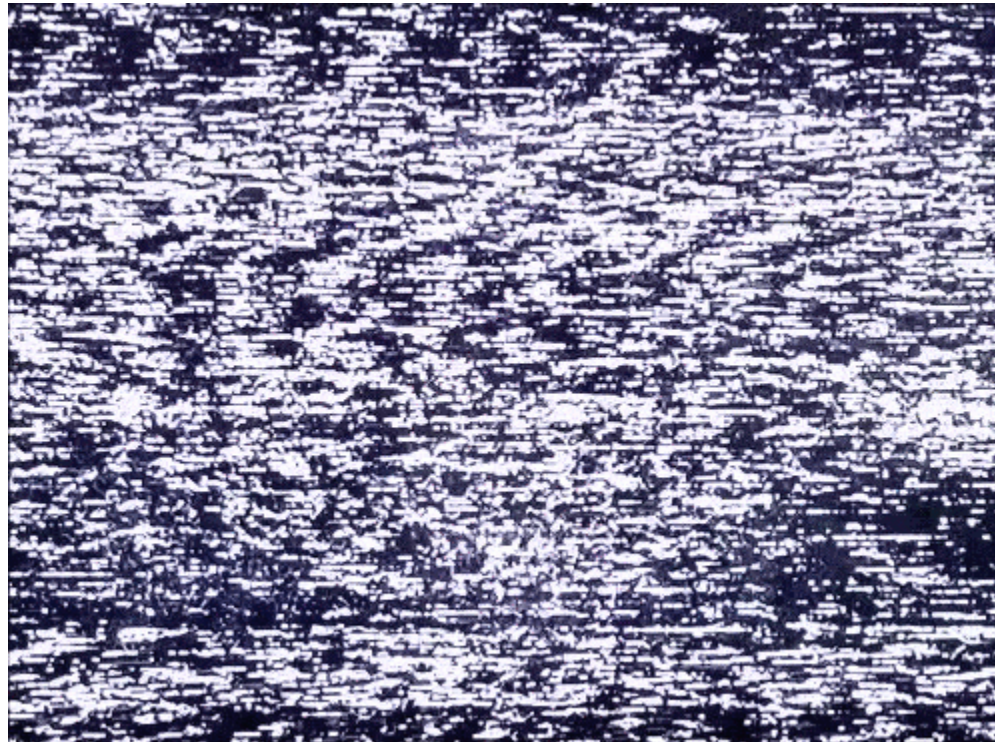
$$R = V \begin{bmatrix} 1 & & \\ & 1 & \\ & & \text{sign}(\det(VU^T)) \end{bmatrix} U^T$$

using  $A = USV^T$  so actually not a *minimal* problem (but certainly a *geometric* one)

- Use 5-point algorithm to get  $[R \ t]$  of the first two cameras
- Keep using P3P for each consecutive camera
- Keep aligning and concatenating



- Use 5-point algorithm to get  $[R \ t]$  of the first two cameras
- Keep using P3P for each consecutive camera
- Keep aligning and concatenating
- But ...



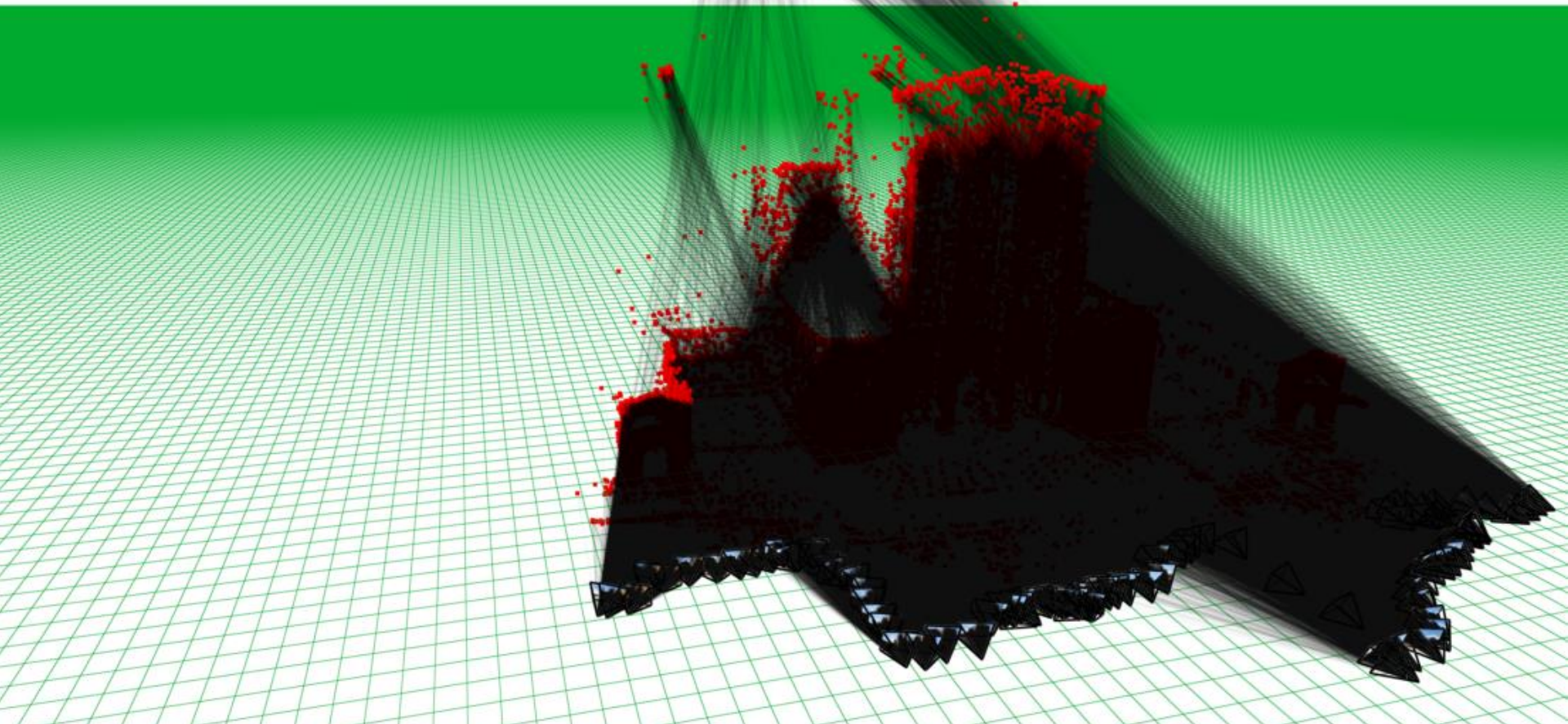
- Use 5-point algorithm to get  $[R \ t]$  of the first two cameras
- Keep using P3P for each consecutive camera
- Keep aligning and concatenating
- Use RANSAC to select the points for the MGPs
- Not enough! Need to do *maximum likelihood estimation* (MLE)

- Bundle Adjustment
  - We have a set of „cameras“
    - Each described by a pose  $[R \ t]$ , projection  $K$
    - Possibly also lens distortion parameters  $l$  and function
  - We have a set of observed points
    - Each described by its position  $X$
  - We have a set of observations
    - Link between point  $X$ , camera  $C$
    - Position of point in the image  $x$
    - Residual to minimize

$$r = \text{err} \left( x, \text{distort} \left( l, \text{dehomog} \left( P \begin{bmatrix} R & t \\ & 1 \end{bmatrix} X \right) \right) \right)$$



- Guildford Cathedral
- 92 poses
- 57957 landmarks



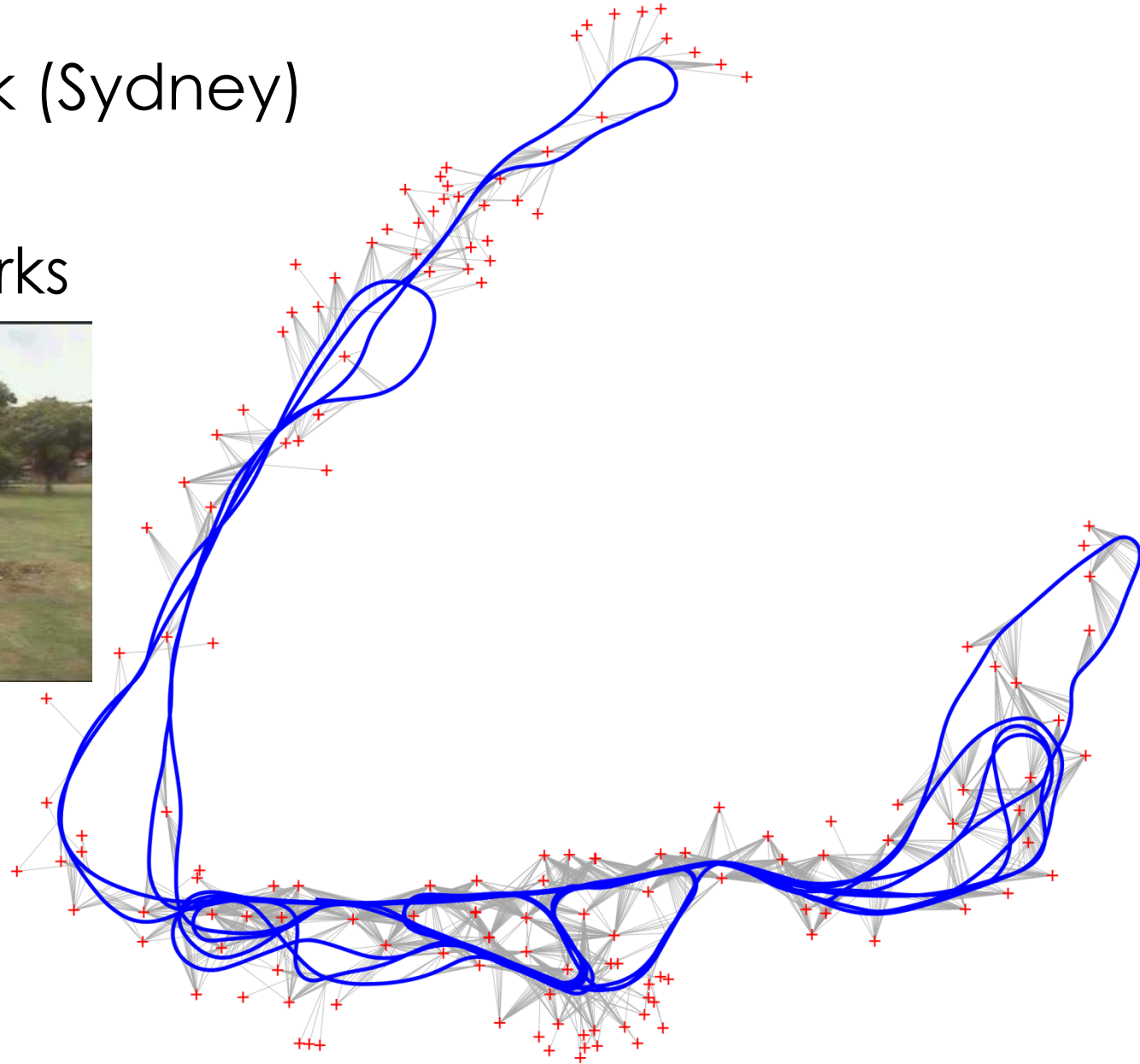
- Simultaneous Localization and Mapping–SLAM
  - We have a set of „robot poses“
    - Each described by a  $[R \ t]$  matrix
  - We *may* have a set of landmarks
    - Each described by its position  $l$
  - We have a set of observations
    - Odometry observations (links between two poses)
      - Estimated distance travelled  $D$  (also a pose itself)
      - Residual to minimize

$$r = \text{err}\left(\begin{bmatrix} R & t \\ & 1 \end{bmatrix}, \begin{bmatrix} R & t \\ & 1 \end{bmatrix} \oplus D\right)$$

- Landmark observations (pose-landmark links)
    - „Measurement“ of the landmark (e.g. Range-bearing vector)
    - Residual to minimize

$$r = \text{err}\left(l, \begin{bmatrix} R & t \\ & 1 \end{bmatrix} \oplus [r \ b]\right)$$

- Victoria Park (Sydney)
- 6969 poses
- 151 landmarks



- In general, inference on *graphical* models
  - Set of variables  $V$ , each contains state vector  $\mathbf{v}_i$
  - Set of edges  $E$ , each a triplet  $\{\mathbf{i}_k, \mathbf{j}_k, \mathbf{o}_k\}$ 
    - Where  $\mathbf{i}$  and  $\mathbf{j}$  are vectors containing vertex indices
  - Minimize

$$\sum_{k=0}^{|E|} \text{err}(\mathbf{v}_{i_k} \ominus \mathbf{v}_{j_k} \oplus \mathbf{o}_k),$$

where  $\ominus$  and  $\oplus$  are vectorial difference and composition operators (e.g. (inverse) matrix multiplication), subject to their precedence and commutativity

- In general, inference on *graphical* models
  - Set of variables  $V$ , each contains state vector  $\mathbf{v}_i$
  - Set of edges  $E$ , each a triplet  $\{\mathbf{i}_k, \mathbf{j}_k, \mathbf{o}_k\}$ 
    - Where  $\mathbf{i}$  and  $\mathbf{j}$  are vectors containing vertex indices
  - Minimize

$$\sum_{k=0}^{|E|} \text{err}(\mathbf{v}_{i_k} \ominus \mathbf{v}_{j_k} \oplus \mathbf{o}_k)$$

- Problem: if there are different kinds of observations, how to make different  $\text{err}(\cdot)$  comparable?
- Solution:  $\text{err}(\mathbf{v}) = \|\mathbf{v}\|_{\Sigma_k} = \mathbf{v} \Sigma_k \mathbf{v}^T$ ,  
where  $\Sigma_k$  is information matrix (inverse covariance)  
for observation  $k$



- Given a set of constraints (observations), find a solution that minimizes L2 error
- For a linear over-determined problem  $A\mathbf{x} = \mathbf{b}$  (where the unknown is  $\mathbf{x}$ )
- The error is  $\mathbf{r} = \mathbf{b} - A\mathbf{x}$ , squared error is\*  
$$\mathbf{r}^2 = (\mathbf{b} - A\mathbf{x})^T (\mathbf{b} - A\mathbf{x}) = \mathbf{b}^T \mathbf{b} - 2\mathbf{x}^T A^T \mathbf{b} + \mathbf{x}^T A^T A \mathbf{x}$$
- This error minimizes where the first derivative cancels  $\mathbf{r}^{2'} = 0$
- Differentiating by  $\mathbf{x}$  gives  $-A^T \mathbf{b} + (A^T A)\mathbf{x} = 0$
- Hence  $\mathbf{x} = (A^T A)^{-1} A^T \mathbf{b} = A^+ \mathbf{b}$
- $A^+$  is the Moore-Penrose pseudo-inverse

\*Note that  $\mathbf{x}^T A^T \mathbf{b}$  is a scalar, thus allowing summation  $\mathbf{x}^T A^T \mathbf{b} + \mathbf{b}^T A \mathbf{x} = 2\mathbf{x}^T A^T \mathbf{b}$ .

- Sometimes we want to weight the observations
- Treated easily, as

$$\mathbf{x} = (A^T W A)^{-1} A^T W \mathbf{b} ,$$

where  $W$  is a diagonal matrix containing the weights

- The weights should be reciprocal variances of the estimated variables (inverse covariances if estimating vectorial quantities)
- Can „bake“ weights into  $A$  by using  $\hat{A} = A\sqrt{W}$  since  $W = W^T$  (diagonal matrix)

- By solving the *normal equation*  $A^T A \mathbf{x} = A^T \mathbf{b}$ 
  - The condition number of  $A^T A$  is greater than of  $A$
  - If  $A$  was sparse but any of its rows is full,  $A^T A$  is dense
- By using orthogonal decompositions
  - Notably SVD of  $A$ ,  $A = USV^T$  and  $A^+ = VS^+U^*$  where  $U^*$  is conjugate transpose (equals  $U^T$  if real)
  - Obtaining  $S^+$  as easy as inverting diagonal entries while skipping zeros
  - So solve  $\mathbf{x} = VS^+U^T \mathbf{b}$ 
    - Slow to compute (expensive Householder reduction to bidiagonal form, followed by iterative diagonalization)
    - May be more precise, can threshold  $S$  to reduce noise
  - Sometimes, there is an additional constraint  $|\mathbf{x}| = 1$ , then only use the smallest singular value, zero the rest

```
O = [0 1; 1 3; 2 3; 3 7];
```

```
% observations of some 1D function
```

```
% O(:,1) are the arguments, a
```

```
% O(:,2) are the desired fit values, b
```

```
% we are trying to estimate  $\mathbf{b} = p + q\mathbf{a} + r\mathbf{a}^2$ 
```

```
% where  $\mathbf{x} = [p \ q \ r]$  is our unknown
```

```
m = length(O);
```

```
A = [ones(m, 1), O(:, 1), O(:, 1).^2];
```

```
% stack the equations  $(1p + \mathbf{a}q + \mathbf{a}^2r)$ 
```

```
 $\mathbf{x} = (A' * A) \setminus A' * O(:,2)$ 
```

```
% solve normal equation (backslash = solve)
```

```
xx = linspace(min(O(:, 1)) - .5, max(O(:, 1)) + .5);
```

```
yy =  $\mathbf{x}(1) + \mathbf{x}(2) * \mathbf{xx} + \mathbf{x}(3) * (\mathbf{xx}.^2)$ ;
```

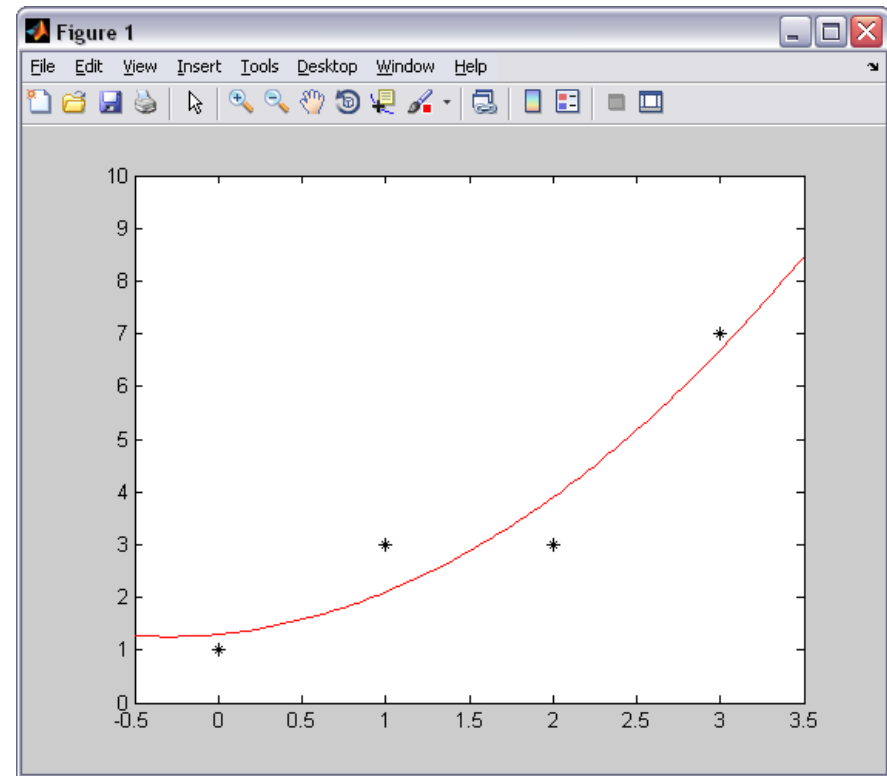
```
% evaluate the estimated model
```

```
plot(O(:, 1), O(:, 2), '*k') % plot a-s, b-s
```

```
hold on
```

```
plot(xx, yy, '-r') % plot the fitted curve
```

```
hold off
```



```
O = [0 1; 1 3; 2 3; 3 7];
```

```
% observations of some 1D function
```

```
% O(:,1) are the arguments, a
```

```
% O(:,2) are the desired fit values, b
```

```
% we are trying to estimate b =  $p + qa + ra^2$ 
```

```
% where x = [p q r] is our unknown
```

```
m = length(O);
```

```
A = [ones(m, 1), O(:, 1), O(:, 1).^2];
```

```
% stack the equations ( $1p + \mathbf{a}q + \mathbf{a}^2r$ )
```

```
[U,S,V] = svd(A); % take SVD of A
```

```
Splus = zeros(size(S')); % S' may be rectangular
```

```
n = size(A, 2); % length of the diagonal
```

```
Splus(1:n, 1:n) = diag(1 ./ diag(S))' % recip. diag.
```

```
x = V * Splus * U' * O(:,2)
```

```
% solve using SVD
```

```
xx = linspace(min(O(:, 1)) - .5, max(O(:, 1)) + .5);
```

```
yy = x(1) + x(2) * xx + x(3) * (xx.^2);
```

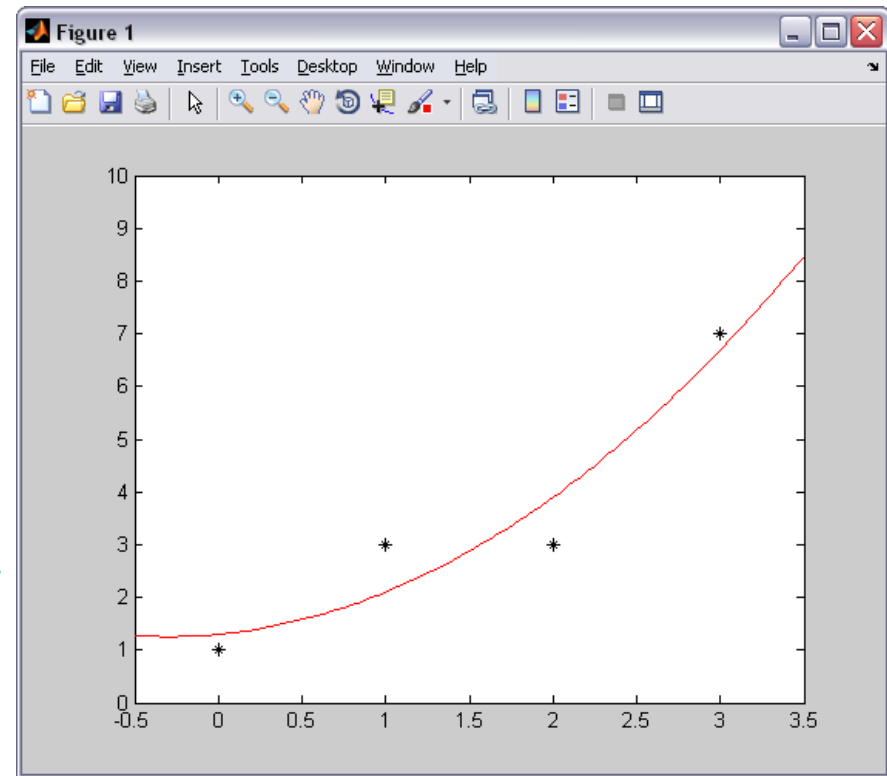
```
% evaluate the estimated model
```

```
plot(O(:, 1), O(:, 2), '*k') % plot a-s, b-s
```

```
hold on
```

```
plot(xx, yy, '-r') % plot the fitted curve
```

```
hold off
```



```
O = [0 1; 1 3; 2 3; 3 7];
```

```
% observations of some 1D function
```

```
% O(:,1) are the arguments, a
```

```
% O(:,2) are the desired fit values, b
```

```
W = diag([10 1 10 10]);
```

```
% weights for the observations
```

```
m = length(O);
```

```
A = [ones(m, 1), O(:, 1), O(:, 1).^2];
```

```
% stack the equations ( $1p + \mathbf{a}q + \mathbf{a}^2r$ )
```

```
x = (A' * W * A) \ A' * W * O(:, 2)
```

```
% solve normal equation (backslash = solve)
```

```
xx = linspace(min(O(:, 1)) - .5, max(O(:, 1)) + .5);
```

```
yy = x(1) + x(2) * xx + x(3) * (xx.^2);
```

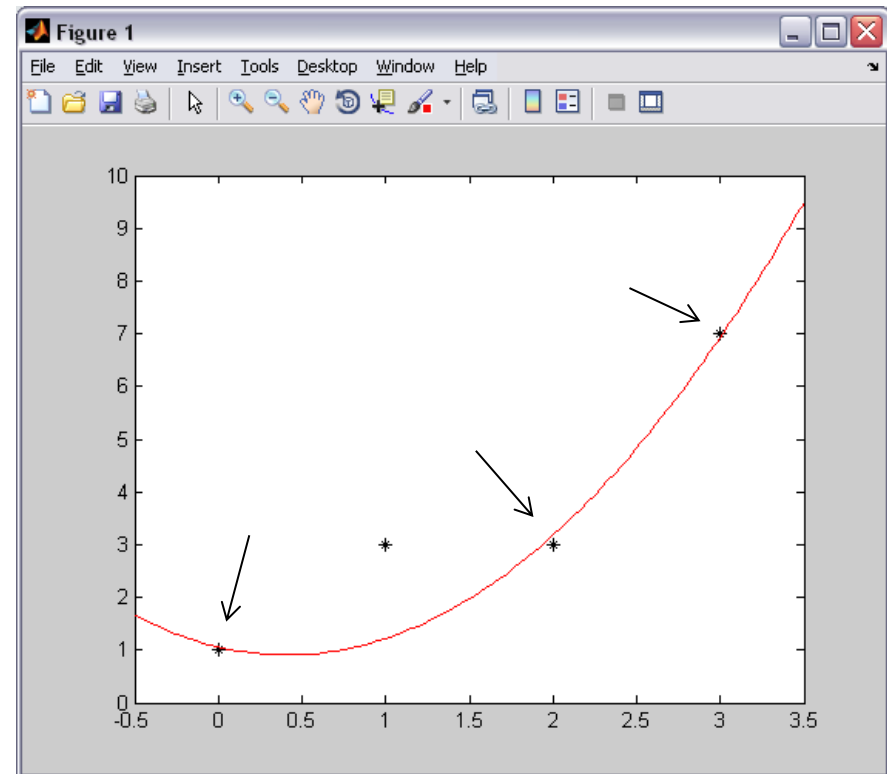
```
% evaluate the estimated model
```

```
plot(O(:, 1), O(:, 2), '*k') % plot a-s, b-s
```

```
hold on
```

```
plot(xx, yy, '-r') % plot the fitted curve
```

```
hold off
```



- The estimation for bundle adjustment is highly *nonlinear* (rotations, projections, ...)
- Formally, instead of linear  $\mathbf{r} = \mathbf{b} - A\mathbf{x}$  we now have  $\mathbf{r} = \mathbf{b} - h(\mathbf{x})$  with nonlinear function  $h(\cdot)$
- To minimize  $s = \mathbf{r}^T \mathbf{r}$ ,  
we again set  $\frac{\partial s}{\partial \mathbf{x}} = 2 \sum \mathbf{r}_i \frac{\partial \mathbf{r}_i}{\partial \mathbf{x}} = 0$
- Instead of directly getting  $\mathbf{x}$ ,  
at iteration  $k$ , we improve  ${}^{k+1}\mathbf{x} = {}^k\mathbf{x} + \Delta\mathbf{x}$
- To solve this, we shall
  - Linearize the problem using Taylor expansion
  - Assume Gaussian noise to be able to do that

- We call  ${}^k\mathbf{x}$  our *linearization point*
- Approximate
$$h(\mathbf{x}) \approx h({}^k\mathbf{x}) + \frac{\partial h({}^k\mathbf{x})}{\partial \mathbf{x}} (\mathbf{x} - {}^k\mathbf{x}) = h({}^k\mathbf{x}) + \mathbf{J}\Delta\mathbf{x}$$
- The Jacobian  $\mathbf{J}$  changes with the linearization
- At step  $k$ , we solve linearized problem
$$\Delta\mathbf{b} = \mathbf{b} - h({}^k\mathbf{x}) ,$$
$$\mathbf{r} = \mathbf{b} - h(\mathbf{x}) = \left( \mathbf{b} - h({}^k\mathbf{x}) \right) + \left( h({}^k\mathbf{x}) - h(\mathbf{x}) \right) ,$$
$$\mathbf{r} \approx \Delta\mathbf{b} - \mathbf{J}\Delta\mathbf{x} \text{ and that is a linear model!}$$
- Solve familiar  $\mathbf{J}^T\mathbf{J}\Delta\mathbf{x} = \mathbf{J}^T\Delta\mathbf{b}$
- Note that also using Hessian  $\mathbf{H} = \frac{\partial^2 h({}^k\mathbf{x})}{\partial \mathbf{x} \partial \mathbf{x}^T}$  would converge faster (this is 1<sup>st</sup> order method)



- Gauss-Newton algorithm

Take initial guess of  ${}^0\mathbf{x}$ , set  $k = 0$

Repeat until end of time

Linearize the system  $\mathbf{J} = \frac{\partial h({}^k\mathbf{x})}{\partial \mathbf{x}}$

Calculate residual  $\Delta\mathbf{b} = \mathbf{b} \ominus h({}^k\mathbf{x})$       note the vectorial op

Solve  $\mathbf{J}^T \mathbf{J} \Delta\mathbf{x} = \mathbf{J}^T \Delta\mathbf{b}$

If  $\|\Delta\mathbf{x}\| < t$  then

stop

${}^{k+1}\mathbf{x} = {}^k\mathbf{x} \oplus \Delta\mathbf{x}$       note the vectorial op

$k = k + 1$

- Unlike linear LS, we need initial guess  ${}^0\mathbf{x}$
- Sometimes, LLS can be used to estimate it
- For BA, there are MGPs to take care of that

```
O = [0 1; 1 3; 2 3; 3 7];
% observations of some 1D function
% O is vector of pairs arg. a and desired fit b

x = [1 -1 -.1]';
% guess 0x (deliberately a bad guess, to take a few steps)
```

% we are trying to estimate  $\mathbf{b} = \mathbf{p} + \mathbf{p}\mathbf{q}\mathbf{a} + \mathbf{p}\mathbf{q}\mathbf{a}^2$   
 % where  $\mathbf{x} = [\mathbf{p} \ \mathbf{q} \ \mathbf{r}]$  is our unknown  
 %  $h(\mathbf{x}) = \mathbf{p} + \mathbf{p}\mathbf{q}\mathbf{a} + \mathbf{p}\mathbf{q}\mathbf{a}^2$ ,  $\mathbf{J} = d\mathbf{h}(\mathbf{x}) / d\mathbf{x} = [1+\mathbf{q}\mathbf{a}+\mathbf{q}\mathbf{r}\mathbf{a}^2, \mathbf{p}\mathbf{a}+\mathbf{p}\mathbf{r}\mathbf{a}^2, \mathbf{q}\mathbf{a}^2]$   
 % we use this parameterization to make it nonlinear (would work with linear  
 % too but would be able to optimize in a single step, which would be boring)

```
plot(O(:, 1), O(:, 2), '*k') % plot a-s, b-s
hold on
xx = linspace(min(O(:, 1)) - .5, max(O(:, 1)) + .5);
for i = 1:10
    yy = x(1) + x(1) * x(2) * xx + x(1) * x(2) * x(3) * (xx.^2);
    plot(xx, yy, '-b') % plot the initial guess
    % evaluate the initial guess

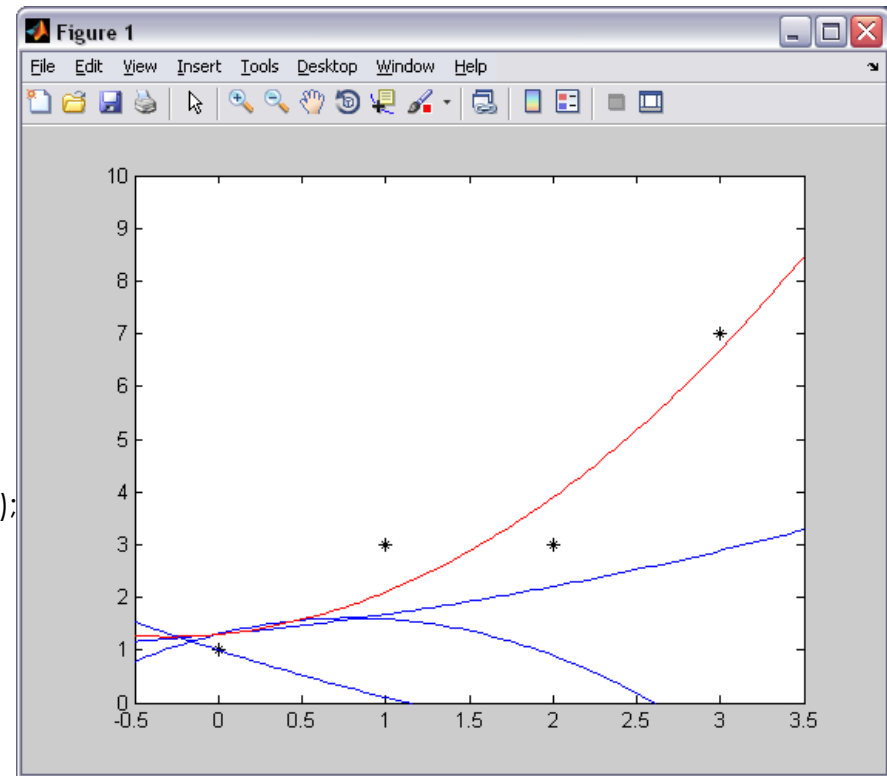
    J = [1 + x(2) * O(:, 1) + x(2) * x(3) * O(:, 1).^2, ...
        x(1) * O(:, 1) + x(1) * x(3) * O(:, 1).^2, x(1) * x(2) * O(:, 1).^2];
    % calculate the Jacobian

    db = O(:,2) - x(1) - O(:,1) * x(1) * x(2) - O(:,1).^2 * x(1) * x(2) * x(3);
    % calculate current error vector

    dx = (J' * J) \ J' * db; norm_dx = norm(dx)
    % solve

    if (norm_dx < 1e-6) % see if we optimize
        break
    end
    x = x + dx % increment
end
```

```
yy = x(1) + x(1) * x(2) * xx + x(1) * x(2) * x(3) * (xx.^2);
plot(xx, yy, '-r') % plot the final in red
hold off
```



And now for something completely different ...

## INTERMEZZO I

- Can readily use Matlab's symbolic toolbox
- Can use a cookbook (e.g. [J. Blanco, 2010, A tutorial on se(3) transformation parameterizations and on-manifold optimization, (TR)].)
- Derivatives know nothing about vectorial, so

$$\frac{\partial f(\mathbf{x})}{\partial \mathbf{x}} := \frac{\partial f(\mathbf{x} \oplus \boldsymbol{\varepsilon})}{\partial \boldsymbol{\varepsilon}}$$

- When dealing with rotations, commutability becomes an issue,  ${}^{k+1}\mathbf{x} = {}^k\mathbf{x} \oplus \Delta\mathbf{x}$  must match  $J$

*„If rotations and translations commuted, we could simply do all our rotations in the morning, before leaving home.“*

- Worked example – pose SLAM
  - We have two poses,  $[R_1 \ t_1]$  and  $[R_2 \ t_2]$
  - We want to avoid optimizing  $4 \times 4$  matrices
  - Internally, they are optimized as  $\mathbb{R}^6$  [axis-angle  $\mathbf{t}$ ]
  - That gives us  $v([R \ \mathbf{t}]) = [\text{aa}(R) \ \mathbf{t}]$  and  $m([\mathbf{r} \ \mathbf{t}]) = [R \ \mathbf{t}]$
  - We have  $\oplus([\mathbf{r}_1 \ \mathbf{t}_1], [\mathbf{r}_2 \ \mathbf{t}_2]) = v(m([\mathbf{r}_1 \ \mathbf{t}_1]) \cdot m([\mathbf{r}_2 \ \mathbf{t}_2]))$   
and  $\ominus([\mathbf{r}_1 \ \mathbf{t}_1], [\mathbf{r}_2 \ \mathbf{t}_2]) = v\left((m([\mathbf{r}_1 \ \mathbf{t}_1]))^{-1} \cdot m([\mathbf{r}_2 \ \mathbf{t}_2])\right)$
  - For NLS, we need  $\mathbf{J} = \frac{\partial h({}^k\mathbf{x})}{\partial \mathbf{x}} = \sum_{k=1}^{|E|} \frac{\partial (x_{i_k} \ominus x_{j_k})}{\partial [x_{i_k} \ x_{j_k}]}$
  - Also, error  $\Delta \mathbf{b} = \sum_{k=1}^{|E|} \mathbf{o}_k \ominus (x_{i_k} \ominus x_{j_k})$
  - Then update  $\hat{\mathbf{x}}_i = \Delta \mathbf{x}_i \oplus \mathbf{x}_i$  for each variable  $x_i \in V$
  - Recall that this is on graph  $\{V, E\}, \{\mathbf{i}_k, \mathbf{j}_k, \mathbf{o}_k\} \in E$

- A simple rule that allows decomposition of derivatives

$$\left(f(g(x))\right)' = f'(g(x)) \cdot g'(x)$$

so instead of calculating long  $\left(f(g(x))\right)'$ , we calculate much shorter  $f'(y)$  and  $g'(x)$  and multiply them numerically when evaluating  $J$

- We typically calculate derivatives of  $\oplus$ ,  $\ominus$ ,  $v(\cdot)$  and  $m(\cdot)$  and chain rule the rest

- Lets try to get a Jacobian of  $v([R \ t]) = [aa(R) \ t]$

```
syms r00 r01 r02 r10 r11 r12 r20 r21 r22 t0 t1 t2 real
```

```
R = [r00 r01 r02; r10 r11 r12; r20 r21 r22]
```

```
t = [t0 t1 t2]'
```

```
% we have a rotation matrix and a translation vector
```

```
qw = sqrt(1 + r00 + r11 + r22) / 2;
```

```
qx = (r21 - r12) / (4 * qw);
```

```
qy = (r02 - r20) / (4 * qw);
```

```
qz = (r10 - r01) / (4 * qw);
```

```
% simplified matrix to quaternion (would need several branches for numerical stability)
```

```
qnorm = simplify(sqrt([qx qy qz] * [qx qy qz]'));;
```

```
halfangle = asin(qnorm);
```

```
% get half of the rotation angle
```

```
ax = simplify(qx / qnorm * 2 * halfangle);
```

```
ay = simplify(qy / qnorm * 2 * halfangle);
```

```
az = simplify(qz / qnorm * 2 * halfangle);
```

```
% get the rotation as axis-angle
```

```
vec = [ax ay az t0 t1 t2]'
```

```
Rtvec = [r00 r01 r02 r10 r11 r12 r20 r21 r22 t0 t1 t2]';
```

```
J = simplify(jacobian(vec, Rtvec));
```

```
% let matlab calculate the derivatives
```

```
ccode(J)
```

```
% export to C
```



[illegible]



- Some more tricks

```
syms qnorm_ qx_ qy_ qz_ qw_ ax_ ay_ az_ real
J = subs(J, ax, ax_)
J = subs(J, ay, ay_)
J = subs(J, az, az_)
J = subs(J, qx, qx_)
J = subs(J, qy, qy_)
J = subs(J, qz, qz_)
J = subs(J, qw, qw_)
J = subs(J, qnorm, qnorm_)
```

% try substituting common subexpressions to tame the beast (don't simplify J in the previous slide!)

```
ccode(J)
% export to C
```

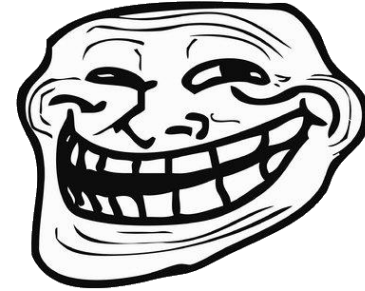
- Still generates 34 kB of C code
- Gotcha – `simple()` and `simplify()` are different
- Lesson learned – don't use  $[R \ t]$ , use Quaternions instead!

42

- Better yet

```
quat_t = [qw qx qy qz t0 t1 t2]
J0 = jacobian(quat_t, Rtvec);
%
syms ax_ ay_ az_ qx_ qy_ qz_ qw_ real
%qnorm = (sqrt([qx qy qz] * [qx qy qz])); % MATLAB trolls you if U no careful ...
qnorm = (sqrt(qx^2 + qy^2 + qz^2));
halfangle_ = asin(qnorm_);
ax_ = (qx_ / qnorm_ * 2 * halfangle_);
ay_ = (qy_ / qnorm_ * 2 * halfangle_);
az_ = (qz_ / qnorm_ * 2 * halfangle_);
vec_ = [ax_ ay_ az_ t0 t1 t2]'
quat_t_ = [qw_ qx_ qy_ qz_ t0 t1 t2]
J1 = jacobian(vec_, quat_t_);
% calculate the jacobians separately (note how J1 is on separate variables)

% chain rule as:
% f( g(R))' = f'( g(R)) * g'(R)
% aa(quat(R))' = aa'(quat(R)) * quat'(R)
% J = J1 * J0
```



```
J1 = subs(J1, qx_, qx); J1 = subs(J1, qy_, qy); J1 = subs(J1, qz_, qz); J1 = subs(J1, qw_, qw)
% substitute the quaternion to J1 (in practice, we would substitute *values* rather than *formulas*)
```

```
error = simplify(J - J1 * J0) % verify, this prints a matrix of zeros
```

- Generates 1.94 and 2.29 kB of C for J0 and J1

```
T[0][0] = -(r21-r12)/sqrt(pow(1.0+r00+r11+r22,3.0))/4.0;  T[0][1] = 0.0;  T[0][2] = 0.0;  T[0][3] = 0.0;  T[0][4] = -(r21-
r12)/sqrt(pow(1.0+r00+r11+r22,3.0))/4.0;  T[0][5] = -1/(sqrt(1.0+r00+r11+r22))/2.0;  T[0][6] = 0.0;  T[0][7] =
1/(sqrt(1.0+r00+r11+r22))/2.0;  T[0][8] = -(r21-r12)/sqrt(pow(1.0+r00+r11+r22,3.0))/4.0;  T[0][9] = 0.0;  T[0][10] = 0.0;  T[0][11]
= 0.0;  T[1][0] = -(r02-r20)/sqrt(pow(1.0+r00+r11+r22,3.0))/4.0;  T[1][1] = 0.0;  T[1][2] = 1/(sqrt(1.0+r00+r11+r22))/2.0;  T[1][3]
= 0.0;  T[1][4] = -(r02-r20)/sqrt(pow(1.0+r00+r11+r22,3.0))/4.0;  T[1][5] = 0.0;  T[1][6] = -1/(sqrt(1.0+r00+r11+r22))/2.0;
T[1][7] = 0.0;  T[1][8] = -(r02-r20)/sqrt(pow(1.0+r00+r11+r22,3.0))/4.0;  T[1][9] = 0.0;  T[1][10] = 0.0;  T[1][11] = 0.0;  T[2][0]
= -(r10-r01)/sqrt(pow(1.0+r00+r11+r22,3.0))/4.0;  T[2][1] = -1/(sqrt(1.0+r00+r11+r22))/2.0;  T[2][2] = 0.0;  T[2][3] =
1/(sqrt(1.0+r00+r11+r22))/2.0;  T[2][4] = -(r10-r01)/sqrt(pow(1.0+r00+r11+r22,3.0))/4.0;  T[2][5] = 0.0;  T[2][6] = 0.0;  T[2][7] =
0.0;  T[2][8] = -(r10-r01)/sqrt(pow(1.0+r00+r11+r22,3.0))/4.0;  T[2][9] = 0.0;  T[2][10] = 0.0;  T[2][11] = 0.0;  T[3][0] = 0.0;
T[3][1] = 0.0;  T[3][2] = 0.0;  T[3][3] = 0.0;  T[3][4] = 0.0;  T[3][5] = 0.0;  T[3][6] = 0.0;  T[3][7] = 0.0;  T[3][8] = 0.0;
T[3][9] = 1.0;  T[3][10] = 0.0;  T[3][11] = 0.0;  T[4][0] = 0.0;  T[4][1] = 0.0;  T[4][2] = 0.0;  T[4][3] = 0.0;  T[4][4] = 0.0;
T[4][5] = 0.0;  T[4][6] = 0.0;  T[4][7] = 0.0;  T[4][8] = 0.0;  T[4][9] = 0.0;  T[4][10] = 1.0;  T[4][11] = 0.0;  T[5][0] = 0.0;
T[5][1] = 0.0;  T[5][2] = 0.0;  T[5][3] = 0.0;  T[5][4] = 0.0;  T[5][5] = 0.0;  T[5][6] = 0.0;  T[5][7] = 0.0;  T[5][8] = 0.0;
T[5][9] = 0.0;  T[5][10] = 0.0;  T[5][11] = 1.0;
```

```
T[0][0] = 2.0/sqrt(qx_*qx_+qy_*qy_+qz_*qz_)*asin(sqrt(qx_*qx_+qy_*qy_+qz_*qz_))-
2.0*qx_*qx_/sqrt(pow(qx_*qx_+qy_*qy_+qz_*qz_,3.0))*asin(sqrt(qx_*qx_+qy_*qy_+qz_*qz_))+2.0*qx_/(qx_*qx_+qy_*qy_+qz_*
qz_)/sqrt(1.0-qx_*qx_-qy_*qy_-qz_*qz_);  T[0][1] = -
2.0*qx_/sqrt(pow(qx_*qx_+qy_*qy_+qz_*qz_,3.0))*asin(sqrt(qx_*qx_~+qy_*qy_+qz_*qz_))*qy_+2.0*qx_/(qx_*qx_+qy_*qy_+qz_*qz
_)*qy_/sqrt(1.0-qx_*qx_~qy_*qy_-qz_*qz_);  T[0][2] = -
2.0*qx_/sqrt(pow(qx_*qx_+qy_*qy_+qz_*qz_,3.0))*asin(sqrt(qx_*qx_~+qy_*qy_+qz_*qz_))*qz_+2.0*qx_/(qx_*qx_+qy_*qy_+qz_*qz
_)*qz_/sqrt(1.0-qx_*qx_~qy_*qy_-qz_*qz_);  T[0][3] = 0.0;  T[0][4] = 0.0;  T[0][5] = 0.0;  T[1][0] = -
2.0*qx_/sqrt(pow(qx_*qx_+qy_*qy_+qz_*qz_,3.0))*asin(sqrt(qx_*qx_~+qy_*qy_+qz_*qz_))*qy_+2.0*qx_/(qx_*qx_+qy_*qy_+qz_*qz
_)*qy_/sqrt(1.0-qx_*qx_~qy_*qy_-qz_*qz_);  T[1][1] = 2.0/sqrt(qx_*qx_+qy_*qy_+qz_*qz_)*asin(sqrt(qx_*qx_+qy_*qy_+qz_*qz_))-
2.0*qy_*qy_/sqrt(pow(qx_*qx_+qy_*qy_+qz_*qz_,3.0))*asin(sqrt(qx_*qx_+qy_*qy_+qz_*qz_))+2.0*qy_/(qx_*qx_+qy_*qy_+qz_*
qz_)/sqrt(1.0-qx_*qx_-qy_*qy_-qz_*qz_);  T[1][2] = -
2.0*qy_/sqrt(pow(qx_*qx_+qy_*qy_+qz_*qz_,3.0))*asin(sqrt(qx_*qx_~+qy_*qy_+qz_*qz_))*qz_+2.0*qy_/(qx_*qx_+qy_*qy_+qz_*qz
_)*qz_/sqrt(1.0-qx_*qx_~qy_*qy_-qz_*qz_);  T[1][3] = 0.0;  T[1][4] = 0.0;  T[1][5] = 0.0;  T[2][0] = -
2.0*qx_/sqrt(pow(qx_*qx_+qy_*qy_+qz_*qz_,3.0))*asin(sqrt(qx_*qx_~+qy_*qy_+qz_*qz_))*qz_+2.0*qx_/(qx_*qx_+qy_*qy_+qz_*qz
_)*qz_/sqrt(1.0-qx_*qx_~qy_*qy_-qz_*qz_);  T[2][1] = -
2.0*qy_/sqrt(pow(qx_*qx_+qy_*qy_+qz_*qz_,3.0))*asin(sqrt(qx_*qx_~+qy_*qy_+qz_*qz_))*qz_+2.0*qy_/(qx_*qx_+qy_*qy_+qz_*qz
_)*qz_/sqrt(1.0-qx_*qx_~qy_*qy_-qz_*qz_);  T[2][2] = 2.0/sqrt(qx_*qx_+qy_*qy_+qz_*qz_)*asin(sqrt(qx_*qx_+qy_*qy_+qz_*qz_))-
2.0*qz_*qz_/sqrt(pow(qx_*qx_+qy_*qy_+qz_*qz_,3.0))*asin(sqrt(qx_*qx_+qy_*qy_+qz_*qz_))+2.0*qz_/(qx_*qx_+qy_*qy_+qz_*
qz_)/sqrt(1.0-qx_*qx_-qy_*qy_-qz_*qz_);  T[2][3] = 0.0;  T[2][4] = 0.0;  T[2][5] = 0.0;  T[3][0] = 0.0;  T[3][1] = 0.0;  T[3][2]
= 0.0;  T[3][3] = 1.0;  T[3][4] = 0.0;  T[3][5] = 0.0;  T[4][0] = 0.0;  T[4][1] = 0.0;  T[4][2] = 0.0;  T[4][3] = 0.0;  T[4][4] =
1.0;  T[4][5] = 0.0;  T[5][0] = 0.0;  T[5][1] = 0.0;  T[5][2] = 0.0;  T[5][3] = 0.0;  T[5][4] = 0.0;  T[5][5] = 1.0;
```

4 kB

```
syms temp0 temp1 temp2 temp3 real
subexpr(J0, 'temp0')
>> temp0 =
>> 1+r00+r11+r22
>> J0 =
>> [-1/4*(r21-r12)/temp0^(3/2), 0, ...
subexpr(ans, 'temp1')
subexpr(ans, 'temp2')
...
J0_ = ans;
ccode(temp0)
ccode(temp1)
ccode(temp2)
ccode(J0_)
```

```
T.setZero();  temp5 = 1.0+r00+r11+r22;  temp6 = sqrt(temp5);  temp7 = pow(temp5, 3.0/2);  
T[0][0] = -(r21-r12)/temp7/4.0; T[0][4] = -(r21-r12)/temp7/4.0;  T[0][5] = -0.5/temp6;  
T[0][7] = 0.5/temp6;  T[0][8] = -(r21-r12)/temp7/4.0; T[1][0] = -(r02-r20)/temp7/4.0;  
T[1][2] = 0.5/temp6;  T[1][4] = -(r02-r20)/temp7/4.0;  T[1][6] = -0.5/temp6;  
T[1][8] = -(r02-r20)/temp7/4.0;  T[2][0] = -(r10-r01)/temp7/4.0;  T[2][1] = -0.5/temp6;  
T[2][3] = 0.5/temp6;  T[2][4] = -(r10-r01)/temp7/4.0;  T[2][8] = -(r10-r01)/temp7/4.0;  
T[3][9] = 1.0;  T[4][10] = 1.0;  T[5][11] = 1.0;
```

```
T.setZero();  temp0 = qx_*qx_+qy_*qy_+qz_*qz;  temp1 = sqrt(temp0);  
temp2 = asin(temp1);  temp3 = sqrt(1.0-temp0);  temp4 = sqrt(pow(temp0,3.0))*temp2;  
T[0][0] = 2.0/temp1*temp2-2.0*qx_*qx_/temp4+2.0*qx_*qx_/(temp0)/temp3;  
T[0][1] = -2.0*qx_/temp4*qy_+2.0*qx_/(temp0)*qy_/temp3;  
T[0][2] = -2.0*qx_/temp4*qz_+2.0*qx_/(temp0)*qz_/temp3;  
T[1][0] = -2.0*qx_/temp4*qy_+2.0*qx_/(temp0)*qy_/temp3;  
T[1][1] = 2.0/temp1*temp2-2.0*qy_*qy_/temp4+2.0*qy_*qy_/(temp0)/temp3;  
T[1][2] = -2.0*qy_/temp4*qz_+2.0*qy_/(temp0)*qz_/temp3;  
T[2][0] = -2.0*qx_/temp4*qz_+2.0*qx_/(temp0)*qz_/temp3;  
T[2][1] = -2.0*qy_/temp4*qz_+2.0*qy_/(temp0)*qz_/temp3;  
T[2][2] = 2.0/temp1*temp2-2.0*qz_*qz_/temp4+2.0*qz_*qz_/(temp0)/temp3;  
T[3][3] = 1.0;  T[4][4] = 1.0;  T[5][5] = 1.0;
```

1 kB



- Much easier for fast prototyping
- Typically not horribly imprecise but can be slow

```
Vector6d x1, x2; // the two variables
```

```
Vector6d expectation = x1  $\ominus$  x2; // the function we're differentiating
```

```
Matrix6d J1; // derivative w.r.t. x1
```

```
for(int i = 0; i < 6; ++ i) {  
    Vector6d eps;  
    eps.setZero(); eps(i) = 1e-9;  
    Vector6d shift1 = eps  $\oplus$  x1; // apply infinitesimal shift (ordering!)  
    Vector6d value = shift1  $\ominus$  x2; // see how that changes the output  
    J1.col(i) = (value - expectation) * 1e+9; // note cwise op  
}  
// do the same for x2
```

- Also good for prototyping but quite slow
- A bit like complex numbers, with special semantics (while  $i^2 = -1$ , we use  $e^2 = 0$ )
- Consider  $f(x) = x^2$ , inject  $y = x + e$
- Evaluate using our “complex” arithmetics  
$$f(y) = (x + e)^2 = x^2 + 2xe + e^2 = \underbrace{x^2}_{\text{value}} + \underbrace{2xe}_{\text{derivative}}$$
- For functions of multiple arguments, we need to put  $e$  in each argument separately and re-evaluate several times
- Rules for multiplication, division, transcendentals (can derive from Taylor series)



- Groups on differentiable manifolds
- Lie group  $SE(3)$  and associated algebra  $\mathfrak{se}(3)$
- $SE(3)$  maps to Euclidean space  $\mathbb{E}^3$  ( $[R \ t]$  are in it)
- $\mathfrak{se}(3)$  is the *tangent space* of real skew-symmetric  $4 \times 4$  matrices
- *Exponential map* goes from  $\mathfrak{se}(3)$  to  $SE(3)$
- *Logarithmic map* goes from  $SE(3)$  to  $\mathfrak{se}(3)$
- *Vectorial operator*\*  $[\mathbf{x}]_{\vee}$  packs skew-sym to vec
- *Cross operator*  $[\mathbf{x}]_{\times}$  goes back to skew-sym
- *Lie bracket*  $[\mathbf{x}, \mathbf{y}] = \mathbf{x}\mathbf{y} - \mathbf{y}\mathbf{x}$  (where  $\mathbf{x}, \mathbf{y} \in \mathfrak{se}(3)$ )

\*several different notations exist

- Let's have skew symmetric  $\mathfrak{so}(3)$  matrix  $A$

$$\mathbf{v} = [a \ b \ c], [\mathbf{v}]_{\times} = A = \begin{bmatrix} 0 & -c & b \\ c & 0 & -a \\ -b & a & 0 \end{bmatrix}$$

- Now, thinking about matrix exponent

$$R = e^A = I + \sum_{i=1} \frac{A^i}{i!} \text{ with } A^0 = I$$

- So exponential map of  $\mathfrak{so}(3)$  yields  $R$ , which is *orthogonal* (i.e. rotation matrix, in  $SO(3)$ )
- There's also a logarithm of matrix

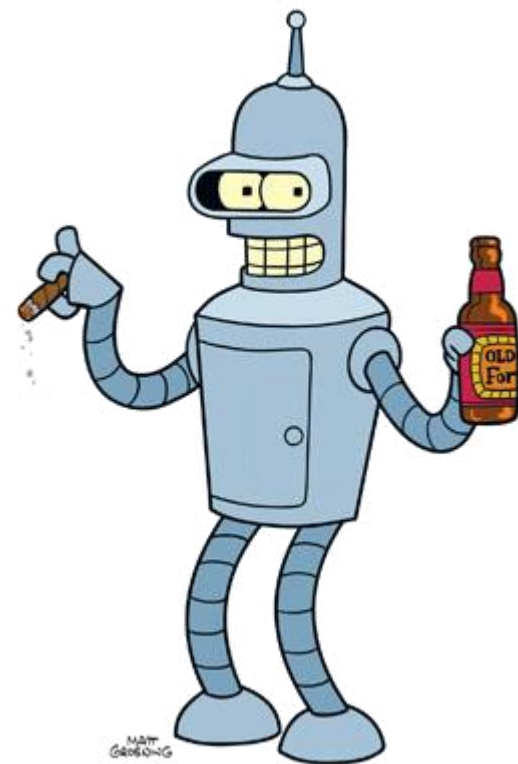
- Following

$$\mathbf{v} = [a \ b \ c], [\mathbf{v}]_{\times} = A = \begin{bmatrix} 0 & -c & b \\ c & 0 & -a \\ -b & a & 0 \end{bmatrix}$$

$$B = \mathbf{v}\mathbf{v}^T = A^2 = \begin{bmatrix} a^2 & ab & ac \\ ab & b^2 & bc \\ ac & bc & c^2 \end{bmatrix}$$

$$e^A = \cos(\theta) I + \frac{\sin(\theta)}{\theta} A + \frac{1 - \cos(\theta)}{\theta^2} B$$

- Seems familiar? :)

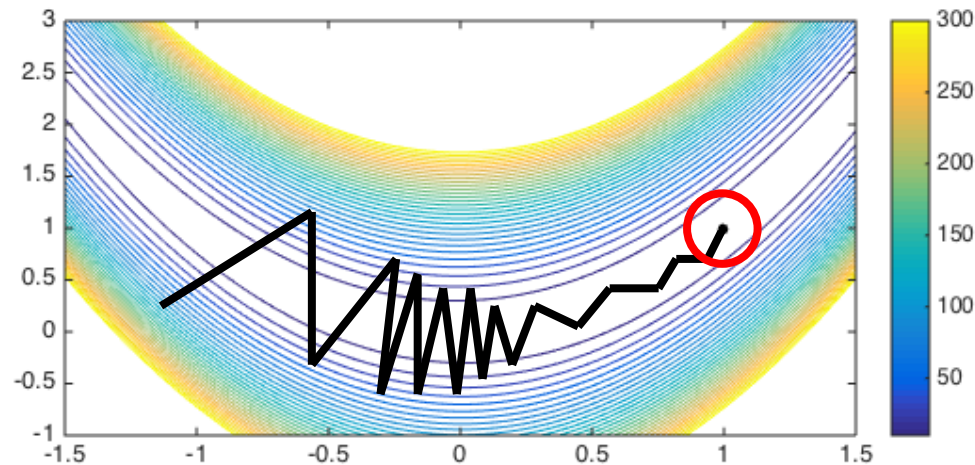


- Going back to derivatives, we can calculate derivatives of  $\exp$ ,  $\log$ , (inverse) compose, just like we did for  $\oplus$ ,  $\ominus$ ,  $v(\cdot)$  and  $m(\cdot)$
- Much of that has nice closed form
- As for least squares, we choose (vectorial)  $se(3)$  or  $sim(3)$  as the internal representation
- See Tom Drummond's TooN library for gritty details and code  
(<https://www.edwardrosten.com/cvd/toon/html-user/>)

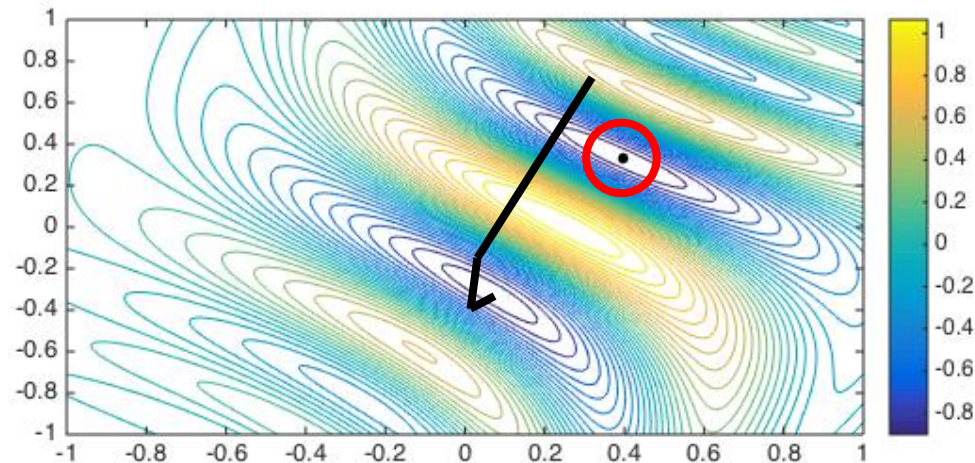
And now ...

# BACK TO BUNDLING

- Gauss-Newton only computes step from *local* gradient, never looks back



Rosenbrock



Chebyshev

This would happen more with gradient descent. This is illustrative.

- Several algorithms that address this
- Levenberg-Marquardt
- Recall ordinary NLS

$$J^T J \Delta \mathbf{x} = J^T \Delta \mathbf{b}$$

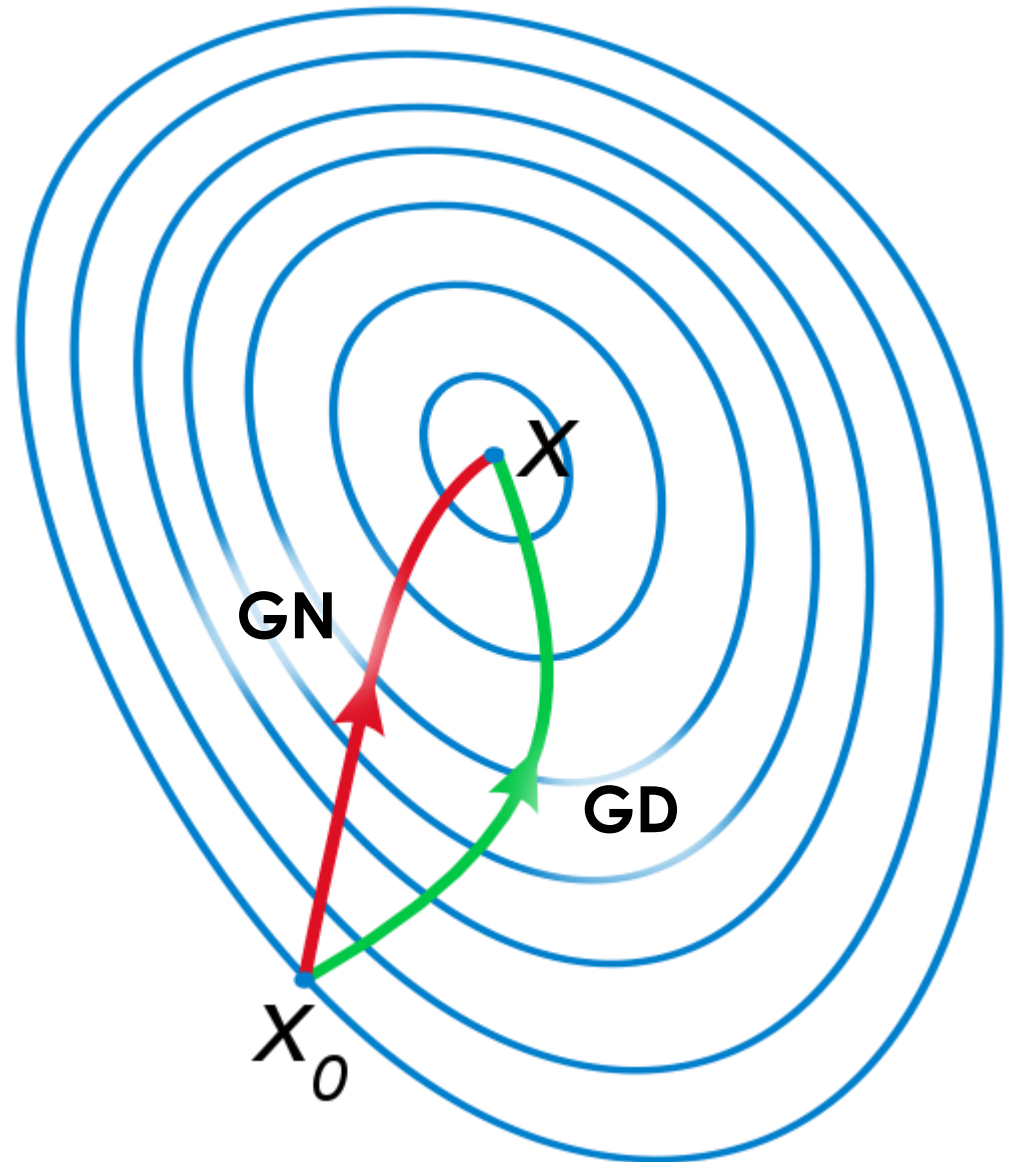
- We're more or less doing  $j^2 x = j b \sim j x = b$
  - We can control step size by magnitude of  $j$
- Levenberg-Marquardt NLS

$$(J^T J + \alpha K) \Delta \mathbf{x} = J^T \Delta \mathbf{b}$$

with common choices  $K = I$  or  $K = \text{diag}(J^T J)$

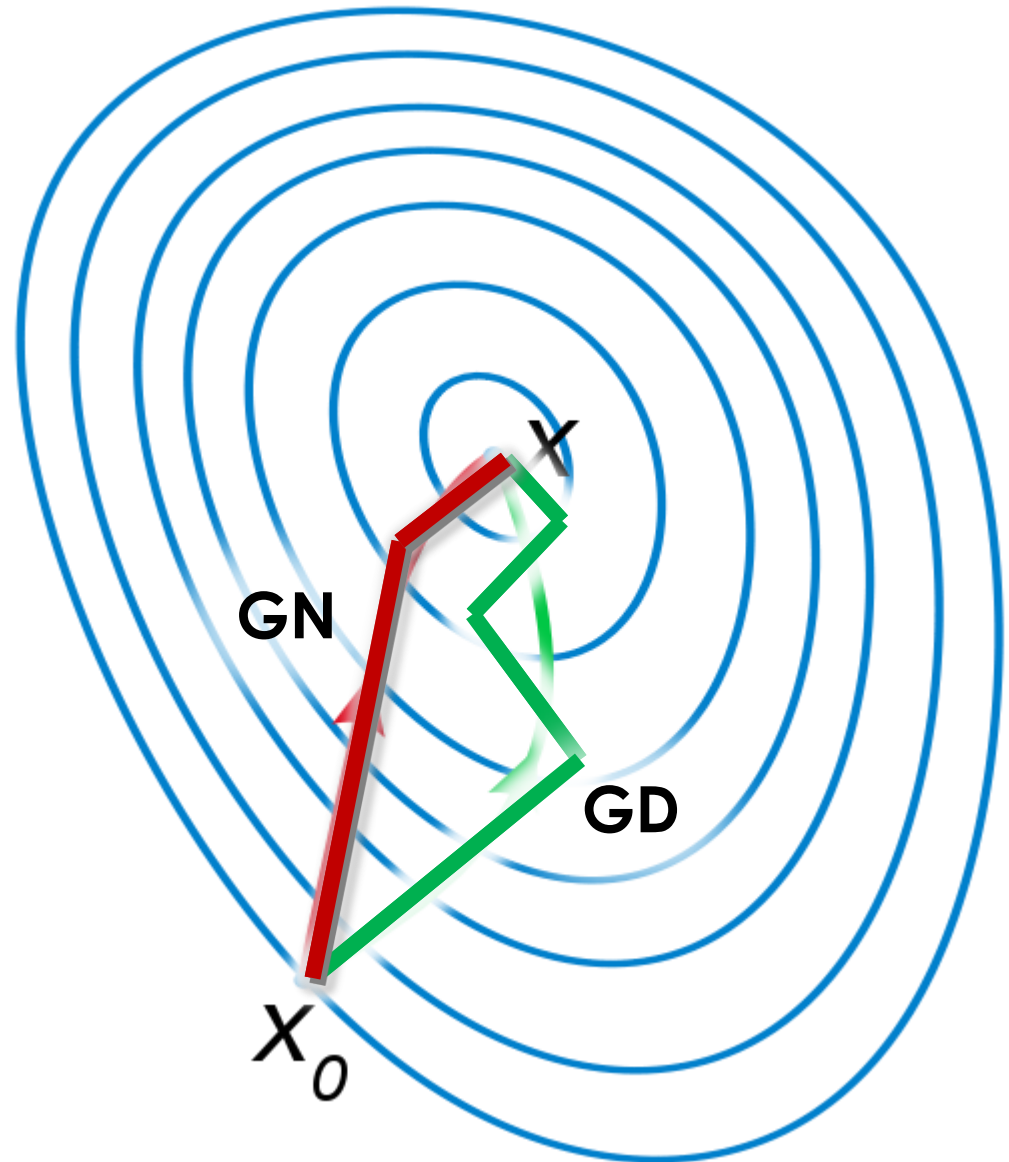
- Changing the value of  $\alpha$  inversely proportionally controls step size and *direction* choosing between GN and steepest descent

- GN typically converges faster
- GD less prone to get stuck but tends to zigzag a lot





- GN typically converges faster
- GD less prone to get stuck but tends to zigzag a lot



# Effects of Levenberg-Marquardt Damping

```
O = [0 1; 1 3; 2 3; 3 7];  
% observations of some 1D function  
% O is vector of pairs arg. a and desired fit b  
  
% we are trying to estimate b =  $p + pqa + pqra^2$   
% where x = [p q r] is our unknown  
%  $h(\mathbf{x}) = p + pqa + pqra^2$ ,  $J = dh(\mathbf{x}) / d\mathbf{x} = [1+qa+qra^2, pa+pra^2, qa^2]$   
% we use this parameterization to make it nonlinear (would work with linear  
% too but would be able to optimize in a single step, which would be boring)
```

```
x = [1 -1 -1]';
```

```
% guess 0x (deliberately a bad guess, to take a few steps)
```

```
plot(O(:, 1), O(:, 2), '*k') % plot a-s, b-s
```

```
hold on
```

```
xx = linspace(min(O(:, 1)) - .5, max(O(:, 1)) + .5);
```

```
for i = 1:100
```

```
yy = x(1) + x(1) * x(2) * xx + x(1) * x(2) * x(3) * (xx.^2);
```

```
plot(xx, yy, '-b') % plot the initial guess
```

```
% evaluate the initial guess
```

```
J = [1 + x(2) * O(:, 1) + x(2) * x(3) * O(:, 1).^2, ...
```

```
      x(1) * O(:, 1) + x(1) * x(3) * O(:, 1).^2, x(1) * x(2) * O(:, 1).^2];
```

```
% calculate the Jacobian
```

```
db = O(:,2) - x(1) - O(:,1) * x(1) * x(2) - O(:,1).^2 * x(1) * x(2) * x(3);
```

```
% calculate current error vector
```

```
dx = (J' * J + 0 * eye(size(J, 2))) \ J' * db; norm_dx = norm(dx)
```

```
% solve
```

```
if (norm_dx < 1e-6) % see if we optimize
```

```
    break
```

```
end
```

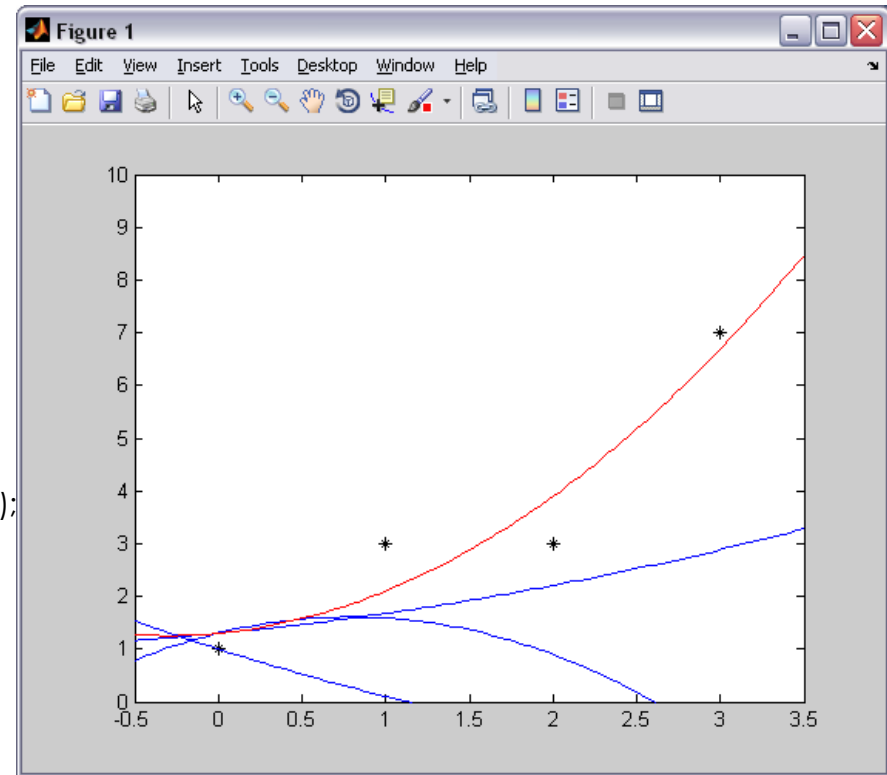
```
x = x + dx % increment
```

```
end
```

```
yy = x(1) + x(1) * x(2) * xx + x(1) * x(2) * x(3) * (xx.^2);
```

```
plot(xx, yy, '-r') % plot the final in red
```

```
hold off
```



# Effects of Levenberg-Marquardt Damping

```
O = [0 1; 1 3; 2 3; 3 7];  
% observations of some 1D function  
% O is vector of pairs arg. a and desired fit b  
  
% we are trying to estimate b =  $p + pqa + pqra^2$   
% where x = [p q r] is our unknown  
%  $h(\mathbf{x}) = p + pqa + pqra^2$ ,  $J = dh(\mathbf{x}) / d\mathbf{x} = [1+qa+qra^2, pa+pra^2, qa^2]$   
% we use this parameterization to make it nonlinear (would work with linear  
% too but would be able to optimize in a single step, which would be boring)
```

```
x = [1 -1 -.1]';
```

```
% guess  $\mathbf{x}$  (deliberately a bad guess, to take a few steps)
```

```
plot(O(:, 1), O(:, 2), '*k') % plot a-s, b-s
```

```
hold on
```

```
xx = linspace(min(O(:, 1)) - .5, max(O(:, 1)) + .5);
```

```
for i = 1:100
```

```
    yy = x(1) + x(1) * x(2) * xx + x(1) * x(2) * x(3) * (xx.^2);
```

```
    plot(xx, yy, '-b') % plot the initial guess
```

```
    % evaluate the initial guess
```

```
    J = [1 + x(2) * O(:, 1) + x(2) * x(3) * O(:, 1).^2, ...
```

```
          x(1) * O(:, 1) + x(1) * x(3) * O(:, 1).^2, x(1) * x(2) * O(:, 1).^2];
```

```
    % calculate the Jacobian
```

```
    db = O(:, 2) - x(1) - O(:, 1) * x(1) * x(2) - O(:, 1).^2 * x(1) * x(2) * x(3);
```

```
    % calculate current error vector
```

```
    dx = (J' * J + 1 * eye(size(J, 2))) \ J' * db; norm_dx = norm(dx)
```

```
    % solve
```

```
    if (norm_dx < 1e-6) % see if we optimize
```

```
        break
```

```
    end
```

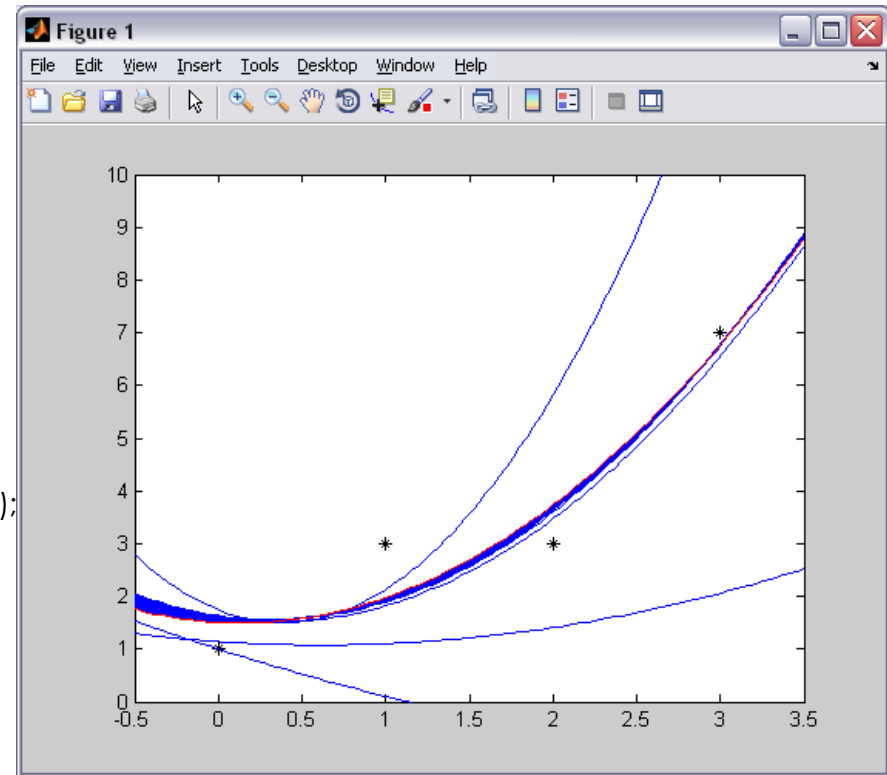
```
    x = x + dx % increment
```

```
end
```

```
yy = x(1) + x(1) * x(2) * xx + x(1) * x(2) * x(3) * (xx.^2);
```

```
plot(xx, yy, '-r') % plot the final in red
```

```
hold off
```



# Effects of Levenberg-Marquardt Damping

```
O = [0 1; 1 3; 2 3; 3 7];  
% observations of some 1D function  
% O is vector of pairs arg. a and desired fit b  
  
% we are trying to estimate b =  $p + pqa + pqra^2$   
% where x = [p q r] is our unknown  
%  $h(\mathbf{x}) = p + pqa + pqra^2$ ,  $J = dh(\mathbf{x}) / d\mathbf{x} = [1+qa+qra^2, pa+pra^2, qa^2]$   
% we use this parameterization to make it nonlinear (would work with linear  
% too but would be able to optimize in a single step, which would be boring)
```

```
x = [1 -1 -.1]';
```

```
% guess 0x (deliberately a bad guess, to take a few steps)
```

```
plot(O(:, 1), O(:, 2), '*k') % plot a-s, b-s
```

```
hold on
```

```
xx = linspace(min(O(:, 1)) - .5, max(O(:, 1)) + .5);
```

```
for i = 1:100
```

```
yy = x(1) + x(1) * x(2) * xx + x(1) * x(2) * x(3) * (xx.^2);
```

```
plot(xx, yy, '-b') % plot the initial guess
```

```
% evaluate the initial guess
```

```
J = [1 + x(2) * O(:, 1) + x(2) * x(3) * O(:, 1).^2, ...
```

```
      x(1) * O(:, 1) + x(1) * x(3) * O(:, 1).^2, x(1) * x(2) * O(:, 1).^2];
```

```
% calculate the Jacobian
```

```
db = O(:,2) - x(1) - O(:,1) * x(1) * x(2) - O(:,1).^2 * x(1) * x(2) * x(3);
```

```
% calculate current error vector
```

```
dx = (J' * J + 10 * eye(size(J, 2))) \ J' * db; norm_dx = norm(dx)
```

```
% solve
```

```
if (norm_dx < 1e-6) % see if we optimize
```

```
    break
```

```
end
```

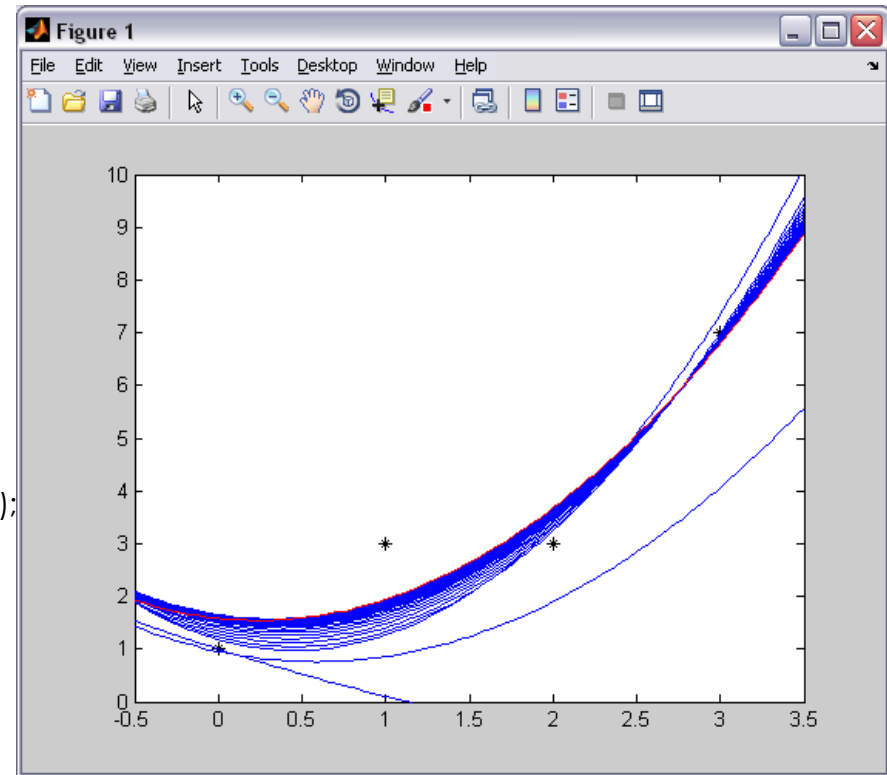
```
x = x + dx % increment
```

```
end
```

```
yy = x(1) + x(1) * x(2) * xx + x(1) * x(2) * x(3) * (xx.^2);
```

```
plot(xx, yy, '-r') % plot the final in red
```

```
hold off
```



# Effects of Levenberg-Marquardt Damping

```
O = [0 1; 1 3; 2 3; 3 7];  
% observations of some 1D function  
% O is vector of pairs arg. a and desired fit b
```

% we are trying to estimate  $\mathbf{b} = \mathbf{p} + \mathbf{pqa} + \mathbf{pqra}^2$   
% where  $\mathbf{x} = [\mathbf{p} \ \mathbf{q} \ \mathbf{r}]$  is our unknown  
%  $h(\mathbf{x}) = \mathbf{p} + \mathbf{pqa} + \mathbf{pqra}^2$ ,  $\mathbf{J} = d\mathbf{h}(\mathbf{x}) / d\mathbf{x} = [1+\mathbf{qa}+\mathbf{qra}^2, \mathbf{pa}+\mathbf{pra}^2, \mathbf{qa}^2]$   
% we use this parameterization to make it nonlinear (would work with linear  
% too but would be able to optimize in a single step, which would be boring)

```
x = [1 -1 -.1]';  
% guess 0x (deliberately a bad guess, to take a few steps)
```

```
plot(O(:, 1), O(:, 2), '*k') % plot a-s, b-s  
hold on
```

```
xx = linspace(min(O(:, 1)) - .5, max(O(:, 1)) + .5);
```

```
for i = 1:100
```

```
    yy = x(1) + x(1) * x(2) * xx + x(1) * x(2) * x(3) * (xx.^2);
```

```
    plot(xx, yy, '-b') % plot the initial guess
```

```
    % evaluate the initial guess
```

```
    J = [1 + x(2) * O(:, 1) + x(2) * x(3) * O(:, 1).^2, ...
```

```
        x(1) * O(:, 1) + x(1) * x(3) * O(:, 1).^2, x(1) * x(2) * O(:, 1).^2];
```

```
    % calculate the Jacobian
```

```
    db = O(:, 2) - x(1) - O(:, 1) * x(1) * x(2) - O(:, 1).^2 * x(1) * x(2) * x(3);
```

```
    % calculate current error vector
```

```
    dx = (J' * J + 100 * eye(size(J, 2))) \ J' * db; norm_dx = norm(dx)
```

```
    % solve
```

```
    if (norm_dx < 1e-6) % see if we optimize
```

```
        break
```

```
    end
```

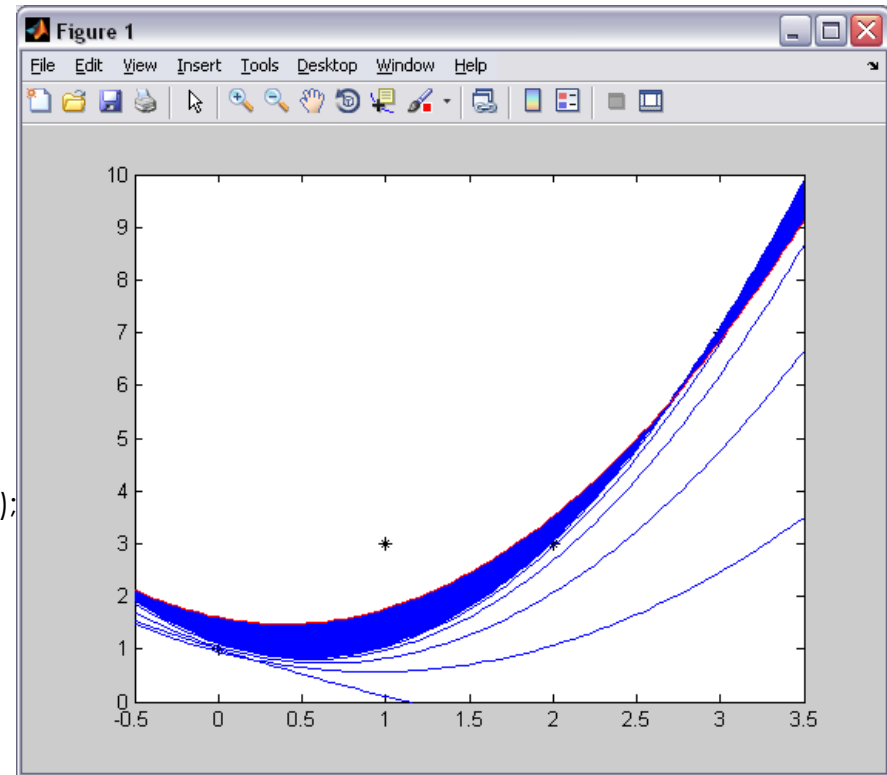
```
    x = x + dx % increment
```

```
end
```

```
yy = x(1) + x(1) * x(2) * xx + x(1) * x(2) * x(3) * (xx.^2);
```

```
plot(xx, yy, '-r') % plot the final in red
```

```
hold off
```



# Effects of Levenberg-Marquardt Damping

```
O = [0 1; 1 3; 2 3; 3 7];  
% observations of some 1D function  
% O is vector of pairs arg. a and desired fit b
```

% we are trying to estimate  $\mathbf{b} = \mathbf{p} + \mathbf{p}\mathbf{q}\mathbf{a} + \mathbf{p}\mathbf{q}\mathbf{a}^2$   
% where  $\mathbf{x} = [\mathbf{p} \ \mathbf{q} \ \mathbf{r}]$  is our unknown  
%  $h(\mathbf{x}) = \mathbf{p} + \mathbf{p}\mathbf{q}\mathbf{a} + \mathbf{p}\mathbf{q}\mathbf{a}^2$ ,  $\mathbf{J} = d\mathbf{h}(\mathbf{x}) / d\mathbf{x} = [1+\mathbf{q}\mathbf{a}+\mathbf{q}\mathbf{r}\mathbf{a}^2, \mathbf{p}\mathbf{a}+\mathbf{p}\mathbf{r}\mathbf{a}^2, \mathbf{q}\mathbf{a}^2]$   
% we use this parameterization to make it nonlinear (would work with linear  
% too but would be able to optimize in a single step, which would be boring)

```
x = [1 -1 -.1]';  
% guess  $\mathbf{x}$  (deliberately a bad guess, to take a few steps)
```

```
plot(O(:, 1), O(:, 2), '*k') % plot a-s, b-s  
hold on
```

```
xx = linspace(min(O(:, 1)) - .5, max(O(:, 1)) + .5);
```

```
for i = 1:100
```

```
yy = x(1) + x(1) * x(2) * xx + x(1) * x(2) * x(3) * (xx.^2);
```

```
plot(xx, yy, '-b') % plot the initial guess
```

```
% evaluate the initial guess
```

```
J = [1 + x(2) * O(:, 1) + x(2) * x(3) * O(:, 1).^2, ...
```

```
      x(1) * O(:, 1) + x(1) * x(3) * O(:, 1).^2, x(1) * x(2) * O(:, 1).^2];
```

```
% calculate the Jacobian
```

```
db = O(:,2) - x(1) - O(:,1) * x(1) * x(2) - O(:,1).^2 * x(1) * x(2) * x(3);
```

```
% calculate current error vector
```

```
dx = (J' * J + 1000 * eye(size(J, 2))) \ J' * db; norm_dx = norm(dx)
```

```
% solve
```

```
if (norm_dx < 1e-6) % see if we optimize
```

```
    break
```

```
end
```

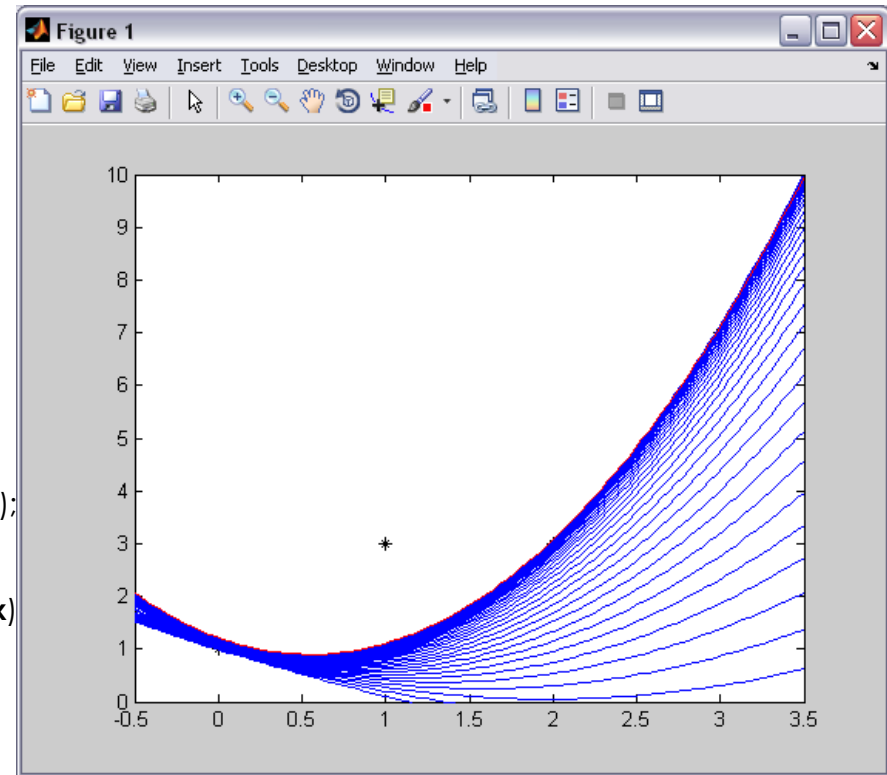
```
x = x + dx % increment
```

```
end
```

```
yy = x(1) + x(1) * x(2) * xx + x(1) * x(2) * x(3) * (xx.^2);
```

```
plot(xx, yy, '-r') % plot the final in red
```

```
hold off
```



```
O = [0 1; 1 3; 2 3; 3 7];
% observations of some 1D function
% O is vector of pairs arg. a and desired fit b
```

% we are trying to estimate  $\mathbf{b} = \mathbf{p} + \mathbf{pqa} + \mathbf{pqra}^2$   
 % where  $\mathbf{x} = [\mathbf{p} \ \mathbf{q} \ \mathbf{r}]$  is our unknown  
 %  $h(\mathbf{x}) = \mathbf{p} + \mathbf{pqa} + \mathbf{pqra}^2$ ,  $\mathbf{J} = d\mathbf{h}(\mathbf{x}) / d\mathbf{x} = [1+\mathbf{qa}+\mathbf{qra}^2, \mathbf{pa}+\mathbf{pra}^2, \mathbf{qa}^2]$   
 % we use this parameterization to make it nonlinear (would work with linear)  
 % too but would be able to optimize in a single step, which would be boring)

```
x = [1 -1 -.1]';
% guess 0x (deliberately a bad guess, to take a few steps)
```

```
plot(O(:, 1), O(:, 2), '*k') % plot a-s, b-s
hold on
```

```
xx = linspace(min(O(:, 1)) - .5, max(O(:, 1)) + .5);
```

```
for i = 1:1000
```

```
    yy = x(1) + x(1) * x(2) * xx + x(1) * x(2) * x(3) * (xx.^2);
```

```
    plot(xx, yy, '-b') % plot the initial guess
```

```
    % evaluate the initial guess
```

```
    J = [1 + x(2) * O(:, 1) + x(2) * x(3) * O(:, 1).^2, ...
```

```
        x(1) * O(:, 1) + x(1) * x(3) * O(:, 1).^2, x(1) * x(2) * O(:, 1).^2];
```

```
    % calculate the Jacobian
```

```
    db = O(:, 2) - x(1) - O(:, 1) * x(1) * x(2) - O(:, 1).^2 * x(1) * x(2) * x(3);
```

```
    % calculate current error vector
```

```
    dx = (J' * J + 1000 * eye(size(J, 2))) \ J' * db; norm_dx = norm(dx)
```

```
    % solve
```

```
    if (norm_dx < 1e-6) % see if we optimize
```

```
        break
```

```
    end
```

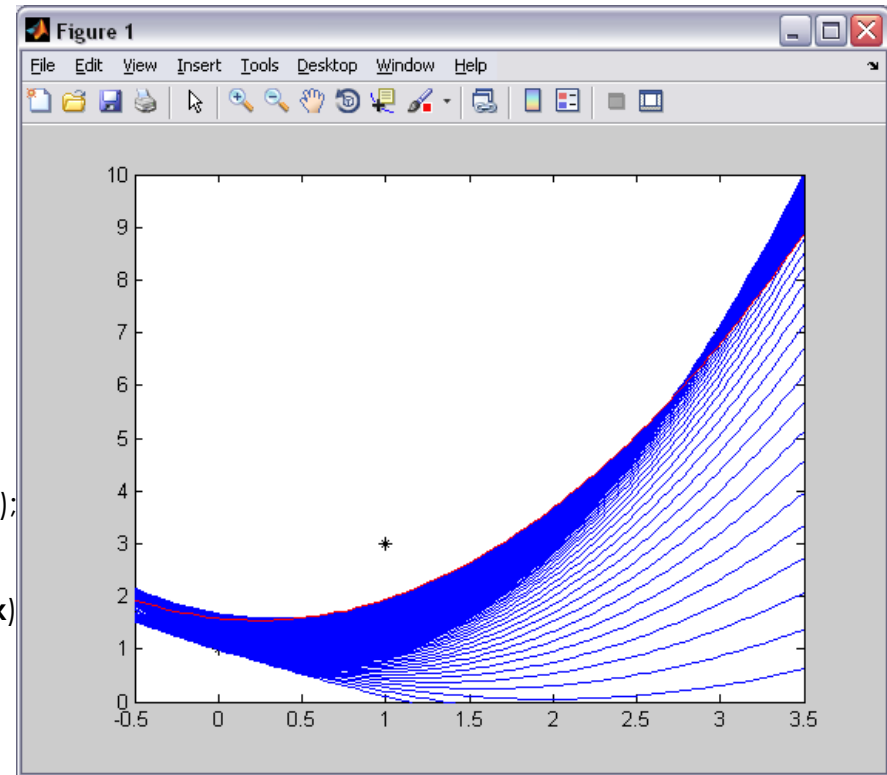
```
    x = x + dx % increment
```

```
end
```

```
yy = x(1) + x(1) * x(2) * xx + x(1) * x(2) * x(3) * (xx.^2);
```

```
plot(xx, yy, '-r') % plot the final in red
```

```
hold off
```



- To control step size, use *optimization gain*

$$g(\Delta \mathbf{x}) = \frac{S(\mathbf{x}) - S(\mathbf{x} + \Delta \mathbf{x})}{L(0) - L(\Delta \mathbf{x})}$$

with linear model error  $L(\Delta \mathbf{x}) = \frac{1}{2} \|h(\mathbf{x}) + J\Delta \mathbf{x}\|^2$   
and squared error  $S(\hat{\mathbf{x}}) = \sum h(\hat{\mathbf{x}})^2$

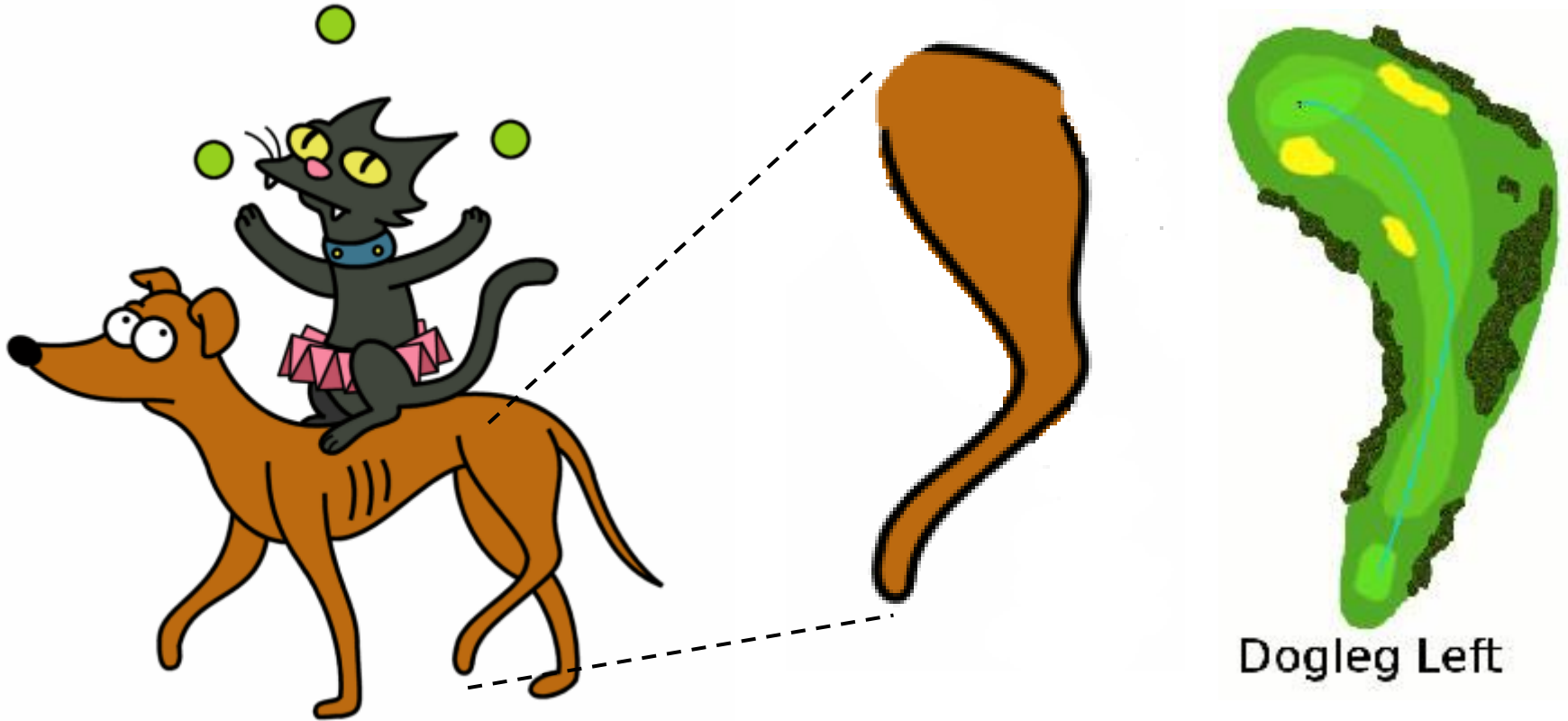
- If the gain is positive, take the step, reduce  $\alpha$
- If the gain is negative, keep the current linearization and increase  $\alpha$  at exponential rate with each failed step
- See [Madsen, 1999, Methods for NLS Problems] for further details





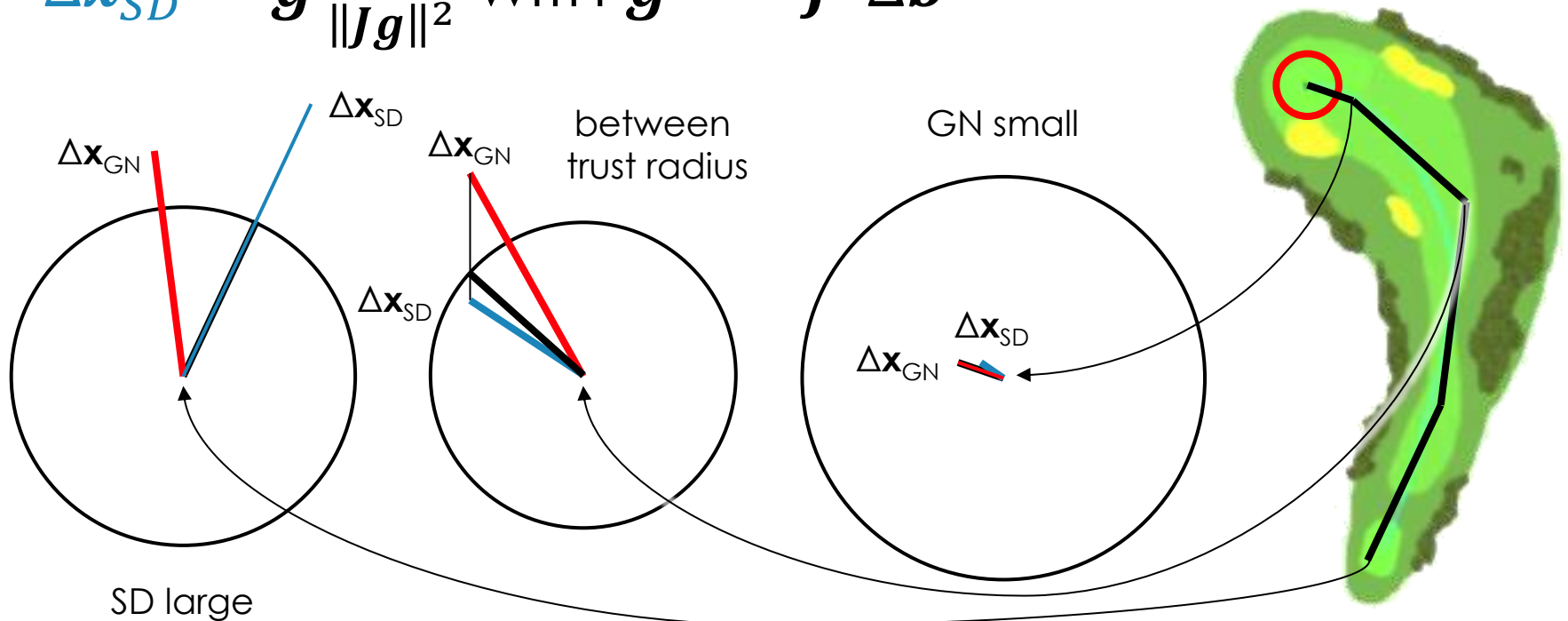
- Quick, **shout** answer to **WIN SAUSAGE!**
- Mike Powell was:
  - A. Cat lady
  - B. Hotdog eating champion
  - C. Golfer
  - D. Computer Scienceman

- Powell's Dogleg (there's also subspace Dogleg)



- Change gradient direction based on magnitude of  $\Delta \mathbf{x}$ , generate step of specified size
- In solving the ordinary normal equation  $\mathbf{J}^T \mathbf{J} \Delta \mathbf{x}_{GN} = \mathbf{J}^T \Delta \mathbf{b}$  we can recover  $\Delta \mathbf{x}_{GN}$  and also

$$\Delta \mathbf{x}_{SD} = \mathbf{g} \frac{\|\mathbf{g}\|^2}{\|\mathbf{J}\mathbf{g}\|^2} \text{ with } \mathbf{g} = -\mathbf{J}^T \Delta \mathbf{b}$$



- Formally, for trust radius  $\underline{\Delta}$ , the step is

$$\Delta \mathbf{x}_{DL} = \begin{cases} \Delta \mathbf{x}_{GN} & \text{if } \|\Delta \mathbf{x}_{GN}\| \leq \underline{\Delta} \\ \underline{\Delta} \cdot \Delta \mathbf{x}_{SD} / \|\Delta \mathbf{x}_{SD}\| & \text{if } \|\Delta \mathbf{x}_{SD}\| \geq \underline{\Delta} \\ \Delta \mathbf{x}_{SD} + \beta (\Delta \mathbf{x}_{GN} - \Delta \mathbf{x}_{SD}) & \text{otherwise} \end{cases}$$

where  $\beta = \frac{-b + \sqrt{b^2 - 4ac}}{2a}$  with  $a = \|\Delta \mathbf{x}_{GN} - \Delta \mathbf{x}_{SD}\|^2$ ,  
 $b = \Delta \mathbf{x}_{SD}^T (\Delta \mathbf{x}_{GN} - \Delta \mathbf{x}_{SD})$  and  $c = \|\Delta \mathbf{x}_{SD}\|^2 - \underline{\Delta}^2$

- $\underline{\Delta}$  is changed based on optimization gain
- The initial trust radius can be large (unlike LM, DL does not resolve the lin. system on bad step)

- Outliers are a problem in computer vision
  - Bad feature matching (mismatched feature)
  - Reflections (matched to a good feature in a mirror)
  - Moving objects (cars, pedestrians, wind)
- Can try to reject suspicious observations
- Can use robust estimation

- Try to identify outliers in the LS framework
- LS minimize *squared* error

$$\mathbf{x} = \operatorname{argmin}_{\mathbf{x}} \frac{1}{2} \sum (b - h(\mathbf{x}))^2$$

big outliers have huge (squared) influence

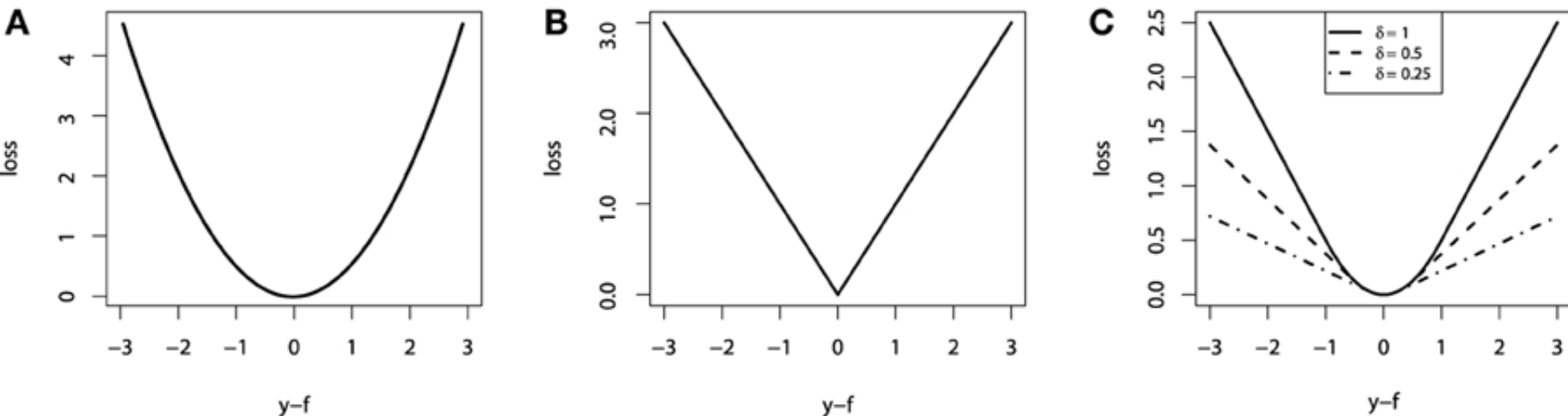
- Try to generalize the error

$$\mathbf{x} = \operatorname{argmin}_{\mathbf{x}} \sum \rho (b - h(\mathbf{x}))$$

where  $\rho(u) = \frac{1}{2}u^2$  is ordinary LS *loss function*

- Try to minimize pseudo-L1 error

$$\rho(u) = \begin{cases} \frac{1}{2}u^2 & \text{if } |u| \leq a \\ \frac{1}{2}a(2|u| - a) & \text{otherwise} \end{cases}$$



Huber kernel

- So how to plug it in? Notice that for LS,  
 $\rho(u) = \frac{1}{2}u^2$  and  $\rho'(u) = u$
- The derivative of loss is *score function*  
 $\psi(u) = \rho'(u)$
- This can be used for weights,

$$\mathbf{J}^T \mathbf{W} \mathbf{J} \Delta \mathbf{x}_{GN} = \mathbf{J}^T \mathbf{W} \Delta \mathbf{b}$$

where  $\mathbf{W} = \text{diag} \left( \frac{\psi(u)}{u} \right)$

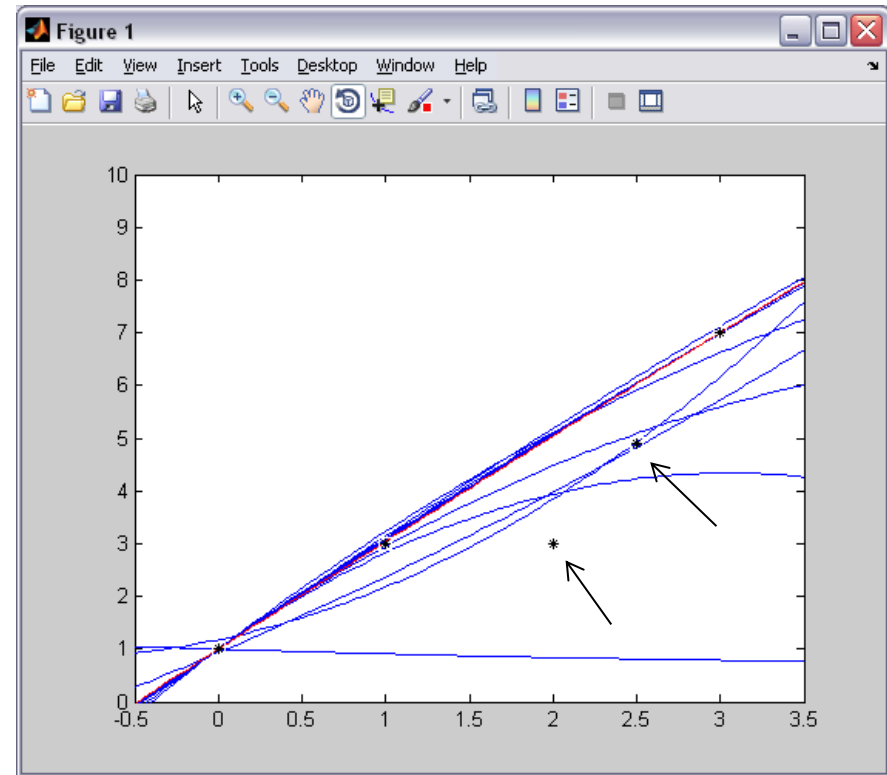
- Let's try setting  $\mathbf{u} = \Delta \mathbf{b}$



```
O = [0 1; 1 3; 2 3; 2.5 4.9; 3 7];  
% observations of some 1D function  
% O(:,1) are the arguments, a  
% O(:,2) are the desired fit values, b
```

```
a = 1.345;  
% guess some parameter for outlier rejection
```

```
for i = 1:10  
    J = ... % calculate the Jacobian  
  
    db = ... % calculate current error vector  
  
    w = zeros(length(db), 1);  
    for j = 1:length(db)  
        u = db(j);  
        if(abs(u) <= a)  
            w(j) = 1 / u; % resolve division by zero  
        else  
            w(j) = a * sign(u) / u;  
        end  
    end  
    W = diag(w); % calculate the robust weights  
  
    dx = (J' * W * J) \ J' * W * db;  
    % solve normal equation (backslash = solve)  
end  
  
% [plot stuff]
```



- Now we're overfitting, the points at the bottom are all *treated as* outliers
- Changing the value of  $a$  is a temporary solution
- We need to approximate *scale of the problem*
- MAD (median absolute deviation)  
 $s = 1.4826 \text{ med}(\text{abs}(\Delta \mathbf{b}))$
- Now set  $\mathbf{u} = \frac{\Delta \mathbf{b}}{s}$
- Alternative to MAD - Huber's second proposal

```
O = [0 1; 1 3; 2 3; 2.5 4.9; 3 7];
```

```
% observations of some 1D function
```

```
% O(:,1) are the arguments, a
```

```
% O(:,2) are the desired fit values, b
```

```
 $\alpha = 1.345;$ 
```

```
% guess some parameter for outlier rejection
```

```
for i = 1:10
```

```
    % [calculate J, db]
```

```
    MAD = median(abs(db));
```

```
    s = 1.4826 * MAD;
```

```
    % calculate MAD
```

```
    w = zeros(length(db), 1);
```

```
    for j = 1:length(db)
```

```
        u = db(j) / s; % apply scale
```

```
        if(abs(u) <=  $\alpha$ )
```

```
            w(j) = 1 / u; % resolve division by zero
```

```
        else
```

```
            w(j) =  $\alpha$  * sign(u) / u;
```

```
        end
```

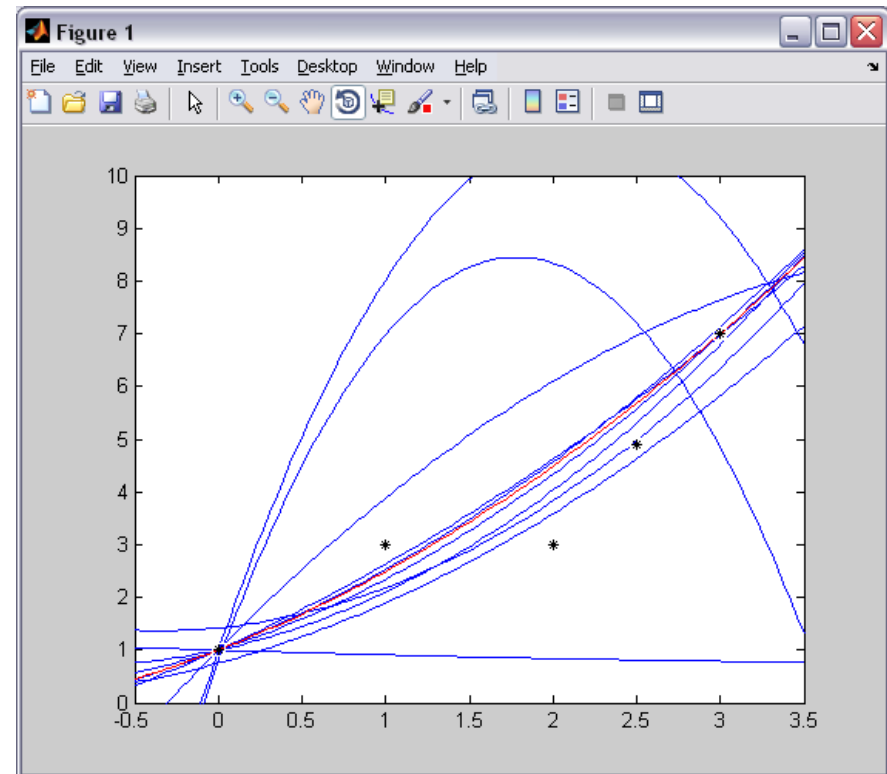
```
    end
```

```
    W = diag(w); % calculate the robust weights
```

```
    % [solve]
```

```
end
```

```
% [plot stuff]
```



```
O = [0 1; 1 7; 2 3; 2.5 4.9; 3 7]; % turn the second point into an OUTLIER
% observations of some 1D function
% O(:,1) are the arguments, a
% O(:,2) are the desired fit values, b
```

```
a = 1.345;
% guess some parameter for outlier rejection
```

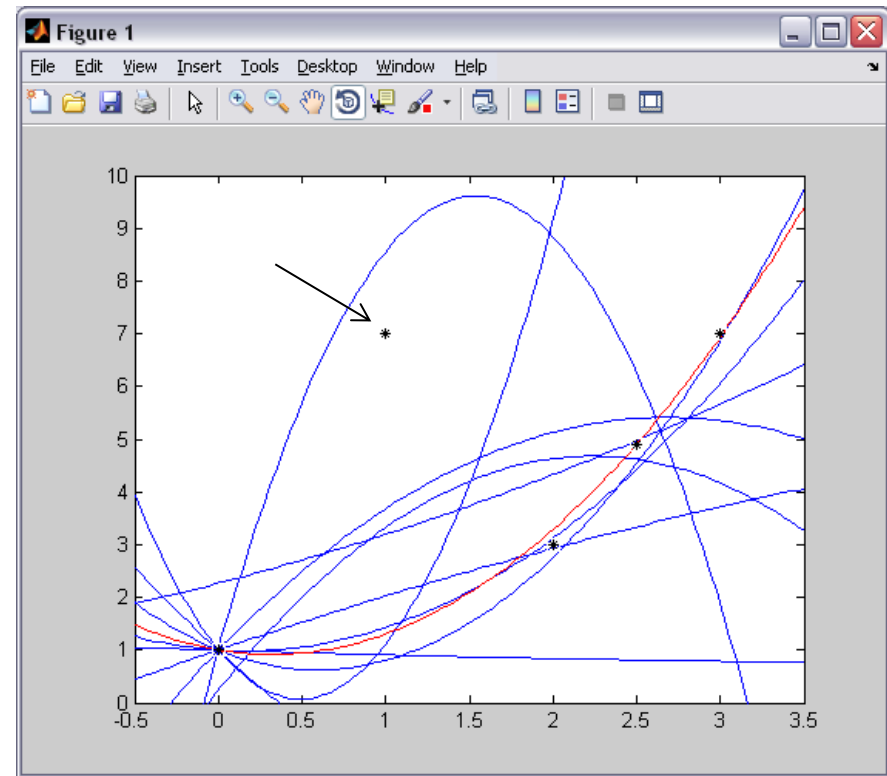
```
for i = 1:10
    % [calculate J, db]

    MAD = median(abs(db));
    s = 1.4826 * MAD;
    % calculate MAD

    w = zeros(length(db), 1);
    for j = 1:length(db)
        u = db(j) / s; % apply scale
        if(abs(u) <= a)
            w(j) = 1 / u; % resolve division by zero
        else
            w(j) = a * sign(u) / u;
        end
    end
    W = diag(w); % calculate the robust weights

    % [solve]
end

% [plot stuff]
```



- What are the magic numbers (1.4826, 1.345)?
- Relative estimator efficiency

$$e = \frac{E((T_2 - \mathbf{x})^2)}{E((T_1 - \mathbf{x})^2)}$$

where  $\mathbf{x}$  is true solution,  $E(\cdot)$  is expectation

- Typically compare to NLS, set efficiency to 95%

- Are there other kernel types?

Table 2.1: A few of the commonly used robust functions. Note that  $a$ ,  $b$  and  $c$  are constant parameters of the individual functions (i.e. not the same variable).

	loss function $\rho(u)$	score function $\psi(u) = \frac{\partial \rho(u)}{\partial u}$
Ordinary LS	$\frac{1}{2}u^2$	$u$
Huber [88]	$\begin{cases} \frac{1}{2}u^2 & \text{if }  u  \leq a \\ \frac{1}{2}a(2 u  - a) & \text{otherwise} \end{cases}$	$\begin{cases} u & \text{if }  u  \leq a \\ a \operatorname{sign}(u) & \text{otherwise} \end{cases}$
Cauchy [83]	$\frac{a^2}{2} \log \left( 1 + \left( \frac{u}{a} \right)^2 \right)$	$\frac{u}{1 + \left( \frac{u}{a} \right)^2}$
Tukey [15]	$\begin{cases} \frac{a^2}{6} \left( 1 - \left( 1 - \left( \frac{u}{a} \right)^2 \right)^3 \right) & \text{if }  u  \leq a \\ \frac{a^2}{6} & \text{otherwise} \end{cases}$	$\begin{cases} u \left( 1 - \left( \frac{u}{a} \right)^2 \right)^2 & \text{if }  u  \leq a \\ 0 & \text{otherwise} \end{cases}$
Hampel [78]	$\begin{cases} \frac{1}{2}u^2 & \text{if }  u  < a \\ a u  - \frac{1}{2}a^2 & \text{if } a \leq  u  < b \\ a \frac{c u  - \frac{1}{2}u^2}{c-b} - \frac{7}{6}a^2 & \text{if } b \leq  u  < c \\ a(b+c-a) & \text{otherwise} \end{cases}$	$\begin{cases} u & \text{if }  u  < a \\ a \operatorname{sign}(u) & \text{if } a \leq  u  < b \\ a \frac{c \operatorname{sign}(u) - u}{c-b} & \text{if } b \leq  u  < c \\ 0 & \text{otherwise} \end{cases}$

And now for something completely different ...

## INTERMEZZO II

- NLS boils down to solving linearized system

$$\mathbf{J}^T \mathbf{J} \Delta \mathbf{x} = \mathbf{J}^T \Delta \mathbf{b}$$

- Most of time is spent there
- In the remainder, let's assume  $\mathbf{A}\mathbf{x} = \mathbf{b}$
- In here,  $\mathbf{A}$  is sparse (most of its entries are zero) but many of the methods apply to dense matrices too



- Different methods
  - Direct methods - reduce the matrix to triangular
    - Proven complexity, mature algorithms
    - Need considerable amounts of memory
  - Iterative methods - do a lot of vector math, iteratively converge to the solution
    - Needs very little memory
    - Convergence depends on problem, preconditioner
    - Quite young field, not much proven
  - Subspace methods - use Eigenvalue-like algorithms to calculate subspace approximations, solve there
- [Davis, 2006, Direct Methods for Sparse Systems]
- [Saad, 2003, Iterative Methods for Sparse Systems]

- Consider the two following systems

$$\begin{bmatrix} 1 & 7 & 5 \\ 1 & 5 & 1 \\ 1 & 2 & 1 \end{bmatrix} \cdot \mathbf{x} = \begin{bmatrix} 9 \\ 6 \\ 9 \end{bmatrix}$$

$$\begin{bmatrix} 1 & 7 & 5 \\ 0 & -2 & -4 \\ 0 & 0 & 6 \end{bmatrix} \cdot \mathbf{x} = \begin{bmatrix} 9 \\ -3 \\ 7.5 \end{bmatrix}$$

- In both cases,  $\mathbf{x} = [9.75, -1, 1.25]$
- What happened?
- Pivoting for stability, cost. Only one r.h.s.

```
A = round(rand(3,3) * 10)
```

```
b = round(rand(3,1) * 10)
```

```
n = length(A);
```

```
for k=1:n-1
```

```
    for i=k+1:n
```

```
        x = A(i,k) / A(k,k);
```

```
        A(i,k+1:n) = A(i,k+1:n) - x * A(k,k+1:n); % row combine
```

```
        A(i,k) = 0; % we just eliminated it
```

```
        b(i) = b(i) - x * b(k);
```

```
    end
```

```
end
```

```
A
```

```
b
```

```
% Gaussian elimination
```

```
x = zeros(n,1);
```

```
for i=n:-1:1 % loop backwards
```

```
    r = b(i) - dot(A(i,i+1:n), x(i+1:n));
```

```
    x(i) = r / A(i, i);
```

```
end
```

```
x
```

```
% back-substitution
```

- Another old method
- Factorize  $\Lambda = LU$ 
  - $L$  being lower triangular with unit diagonal
  - $U$  being upper triangular
- To solve  $\Lambda \mathbf{x} = LU\mathbf{x} = \mathbf{b}$  (and so  $U\mathbf{x} = L^{-1}\mathbf{b}$ )
  - First solve  $L\mathbf{y} = \mathbf{b}$  forward substitution w/o div
  - Then solve  $U\mathbf{x} = \mathbf{y}$  back-substitution
- To work,  $\Lambda$  must be square, invertible
- Will require pivoting for stability
  - Partial  $P\Lambda = LU$  (reorder rows)
  - Full  $P\Lambda Q = LU$  (reorder rows, cols)

```
A = round(rand(3,3) * 10)
```

```
LU = A; % works in-place
```

```
for k = 1:n
```

```
    LU(k+1:n, k) = LU(k+1:n, k) / LU(k, k);
```

```
    % divide by pivot (modifies the rest of this column)
```

```
    for i = k+1:n
```

```
        x = LU(i, k);
```

```
        for j = k+1:n % explicit loop intended
```

```
            LU(i, j) = LU(i, j) - x * LU(k, j); % causes fill-in in sparse version
```

```
        end
```

```
        % reduce the rest of the matrix
```

```
    end
```

```
    % modify the lower-right submatrix
```

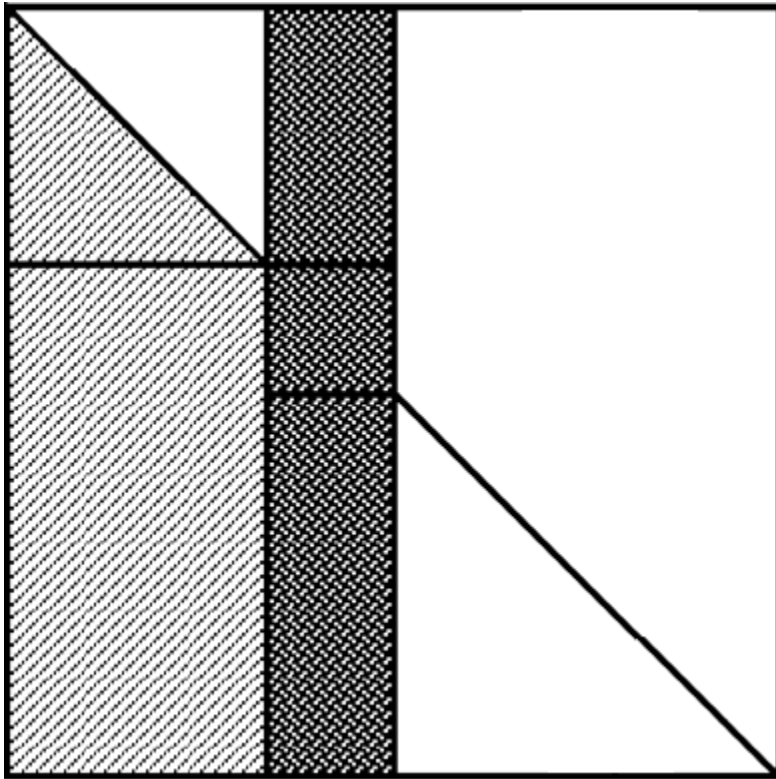
```
end
```

```
U = triu(LU)
```

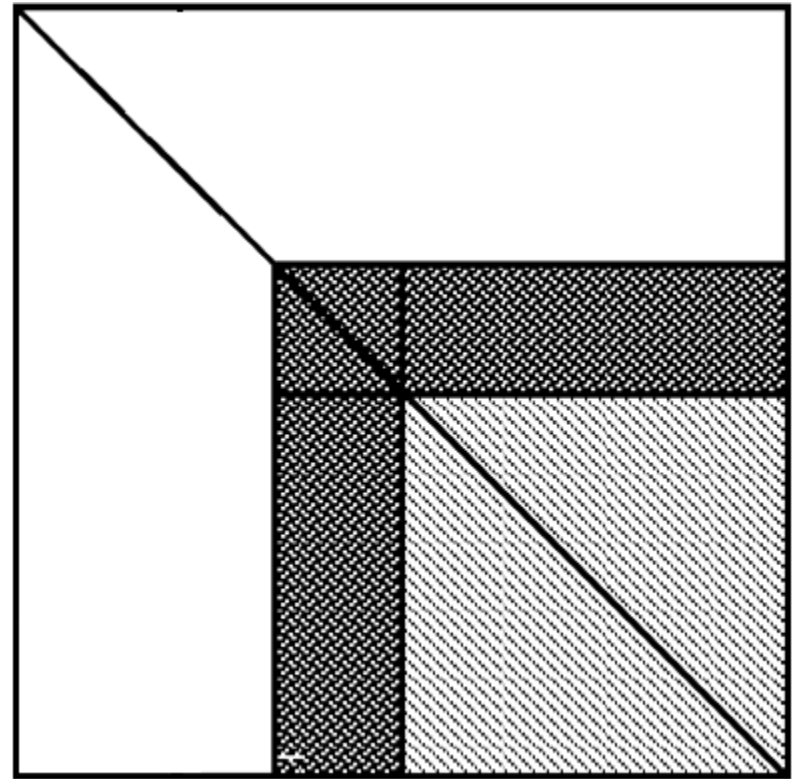
```
L = tril(LU, -1) + eye(size(LU))
```

```
% unpack to L and U, add identity diagonal to L
```

- Several variants, based on ordering of loops **i j k**



Left-looking variant

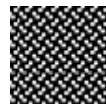


Right-looking variant

**kji**



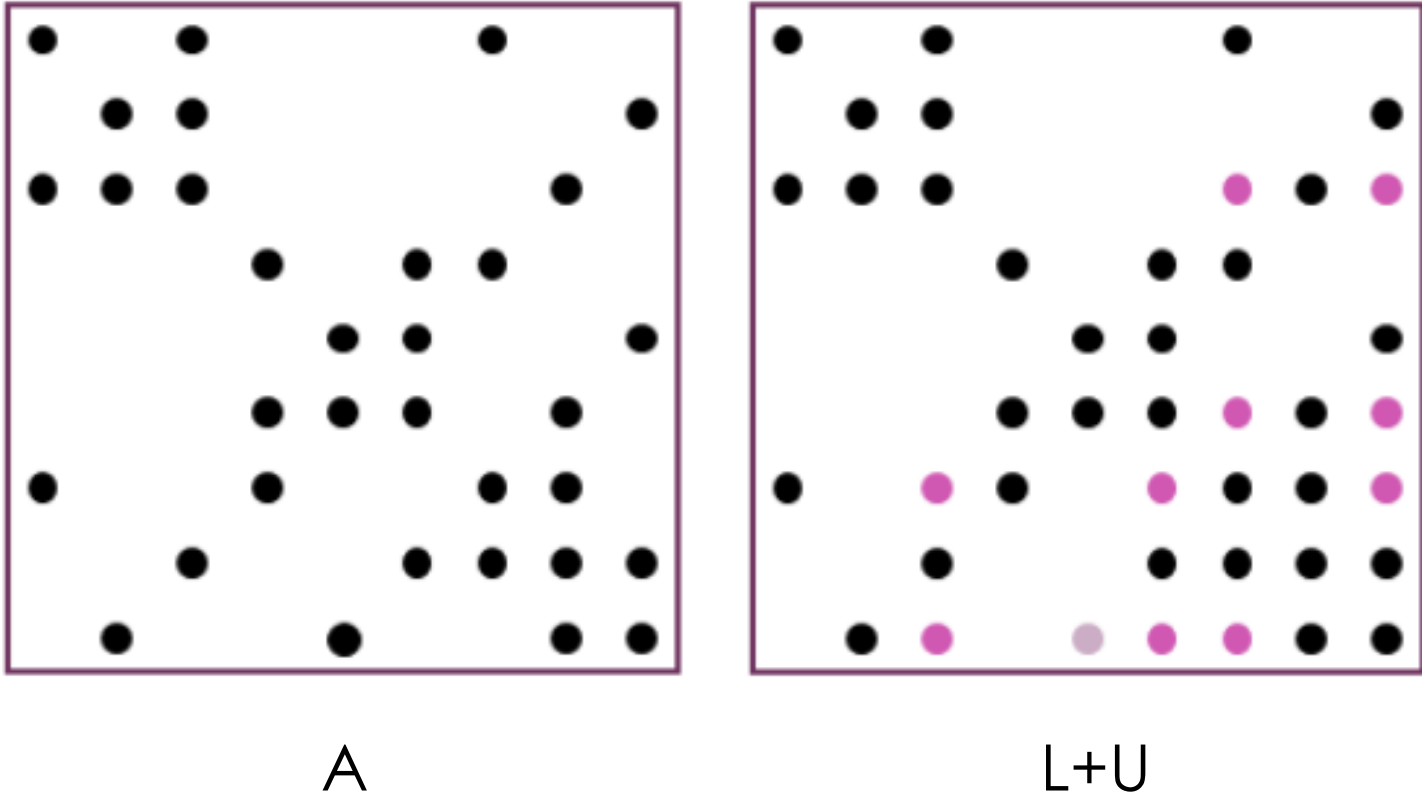
read



write

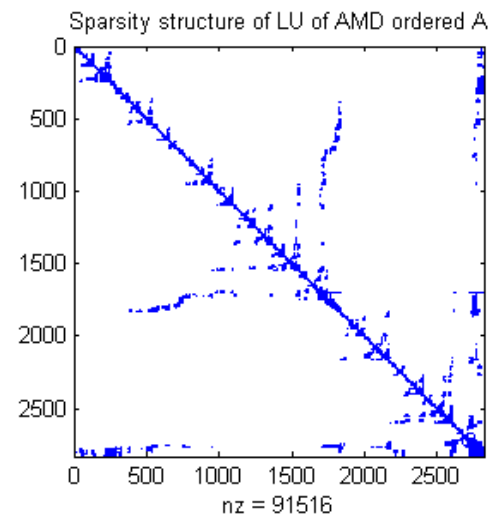
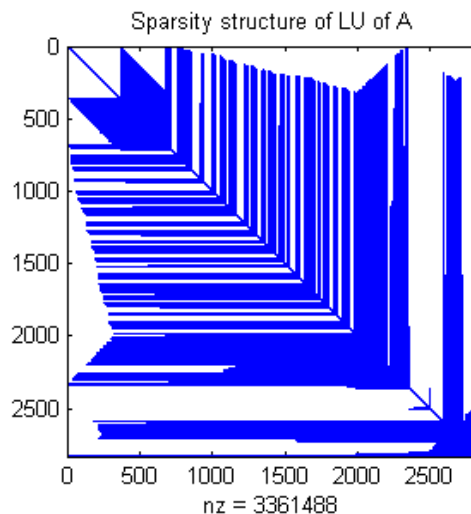
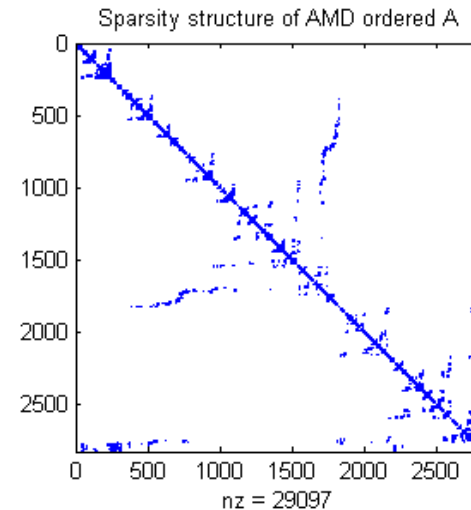
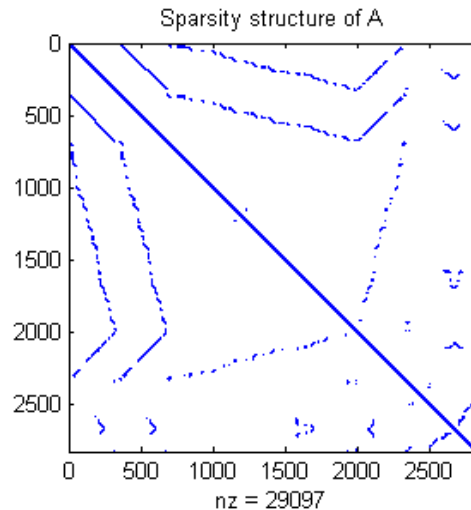
**kij**

- In sparse version, we care about fill-in



- Depends on order of rows / columns

- Approximate Minimum Degree [Amestoy 2004]





enter the characters  
for the symbols shown  
in the box below:



♟=4 ♖=j ♀=d □=n 🎵=x + =k  
▶=v 😊=r ♖=h ♦=t ♣=7 ♗=a

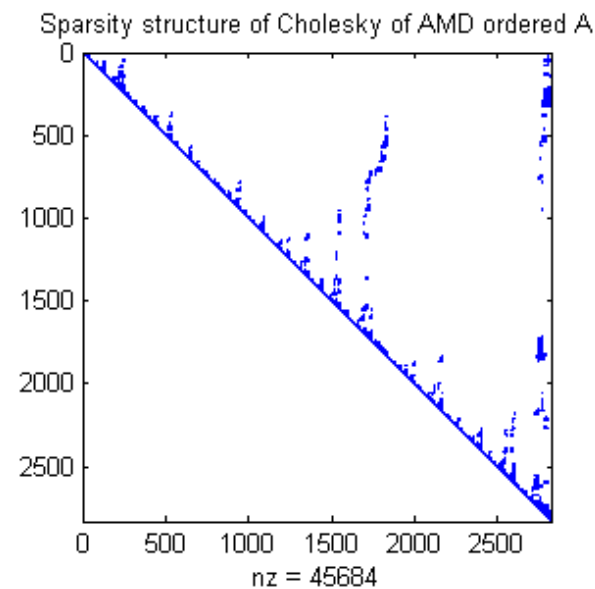
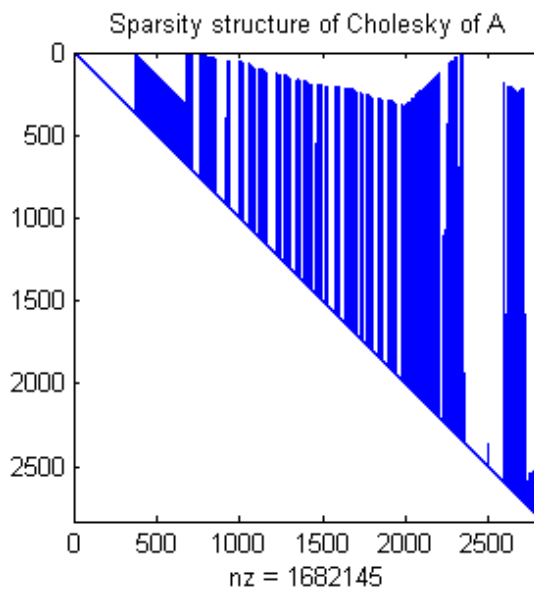
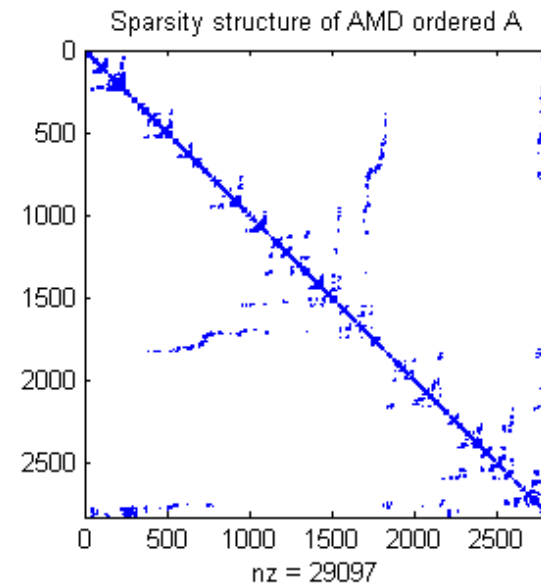
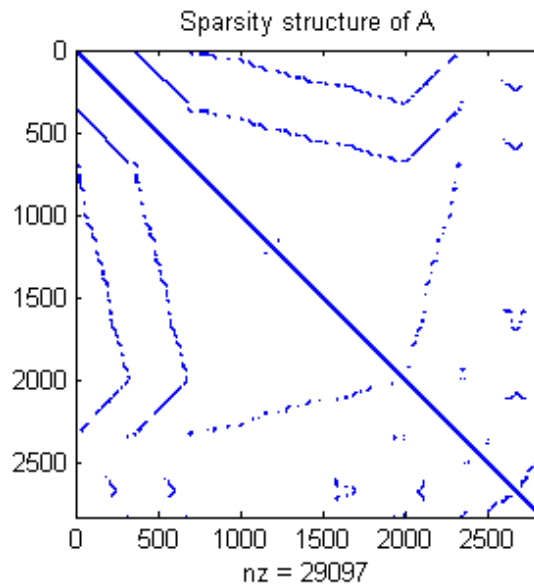
# ENTER CAPTCHA IF YOU'RE AWAKE

- Factorize  $\Lambda = R^T R$ 
  - $R$  being upper triangular, with positive diagonal entries
- To solve  $\Lambda \mathbf{x} = R^T R \mathbf{x} = \mathbf{b}$ 
  - First solve  $R^T \mathbf{y} = \mathbf{b}$  forward substitution
  - Then solve  $R \mathbf{x} = \mathbf{y}$  back-substitution
- To work,  $\Lambda$  must be square, symmetric and positive-definite (SPD)
  - NLS matrices are (up to numerical precision)
  - Adding small damping on the diagonal usually helps
  - Modified Cholesky factorization
- Does *not* need pivoting

```
A = round(rand(3,3) * 10);  
A = A' * A + eye(size(A)) * 10 % try to make it SPD  
  
R = zeros(size(A));  
for j = 1:n % for every column  
    for k = 1:j-1 % for all prev cols that are nnz at row j (know those from etree)  
        s = 0;  
        for i = 1:k-1 % for all blocks above the diagonal in the prev column  
            s = s + R(i, k) * R(i, j); % takes elements from two different columns  
        end  
        % cmod; causes fill-in in the current column  
  
        R(k, j) = (A(k, j) - s) / R(k, k); % accesses upper diagonal of A  
    end  
  
    s = R(1:j-1, j)' * R(1:j-1, j);  
    R(j, j) = sqrt(A(j, j) - s); % must be positive (or modified Cholesky)  
    % cdiv  
end
```

- Elimination tree theory
  - Calculate tree that drives loops in the factorization
  - Saves time
  - Particular orderings also yield elimination trees with independent and balanced subtrees (parallelism)
- Supernodal Cholesky
  - Sometimes, several columns in the factorization have identical nonzero structure - supernodes
  - Treating supernodes as dense allows parallelism, GPU acceleration
- Multifrontal Cholesky
  - Frontal matrices, conceptually similar to supernodes

# Cholesky factorization - AMD again



- Relatively new method (in the 60's)
- Factorize  $A = QR$ 
  - $A$  can now be *rectangular*
  - $Q$  being orthogonal ( $Q^T = Q^{-1}$ )
  - $R$  being upper triangular (zero rows at the bottom)
- To solve  $A\mathbf{x} = QR\mathbf{x} = \mathbf{b}$ 
  - First solve  $\mathbf{y} = Q^T \mathbf{b}$       multiplication\*
  - Then solve  $R\mathbf{x} = \mathbf{y}$       back-substitution
- No pivoting needed
- Fill-in depends on *column* ordering only
  - Out-of-core methods

- Similar to Gaussian elimination
- Householder reflections
  - Compute a reflection about a plane (mirror matrix) that zeroes a lower part of a column in  $R$
  - Minor tricks (sign choice) for numerical stability
  - Record reflections rather than representing  $Q$
- Givens rotations
  - Compute rotation matrix that zeros one element in  $R$
  - Record rotation chains rather than representing  $Q$

- Householder reflections





- *Householder transformation* has form

$$H = I - 2 \frac{vv^T}{v^T v}$$

where  $v$  is nonzero vector

- From definition,  $H = H^T = H^{-1}$ , so  $H$  is both orthogonal and symmetric
- For given vector  $a$ , choose  $v$  so that

$$Ha = \begin{bmatrix} \alpha \\ 0 \\ \vdots \\ 0 \end{bmatrix} = \alpha \begin{bmatrix} 1 \\ 0 \\ \vdots \\ 0 \end{bmatrix} = \alpha e_1$$

- *Givens rotation* operates on pair of rows to introduce single zero
- For given 2-vector  $\mathbf{a} = [a_1 \ a_2]^T$ , if

$$c = \frac{a_1}{\sqrt{a_1^2 + a_2^2}}, \quad s = \frac{a_2}{\sqrt{a_1^2 + a_2^2}}$$

then

$$\mathbf{G}\mathbf{a} = \begin{bmatrix} c & s \\ -s & c \end{bmatrix} \begin{bmatrix} a_1 \\ a_2 \end{bmatrix} = \begin{bmatrix} \alpha \\ 0 \end{bmatrix}$$

- Scalars  $c$  and  $s$  are cosine and sine of angle of rotation, and  $c^2 + s^2 = 1$ , so  $\mathbf{G}$  is orthogonal

- Related to Cholesky
  - $A = QR$
  - $\Lambda = A^T A$
  - $\Lambda = R^T Q^T QR = R^T R$
  - $R$  is the same as in Cholesky factorization, up to the sign of the rows (Cholesky always has positive diag)
- Can directly solve NLS on  $J$  without forming  $J^T J$ 
  - $J = QR$
  - $R\Delta x = Q^T \Delta b$
  - Numerical benefits

```
M = round(rand(5, 3) * 10)

[m n] = size(M);
Q = eye(m, m); % crime against QR: explicit Q
R = M; % works inplace
elim = min(m - 1, n); % decide how many cols to eliminate to get triangular R
for i = 1:elim
    Aii = R(i, i);
    Ai_norm = sqrt(R(i:m, i)' * R(i:m, i));
    dii = abs(Ai_norm) * sign(sign(Aii) + 0.5); % !!! need zero-avoiding sign function !!!
    wii = Aii - dii;
    two_fi_inv_squared = -1 / (wii * dii);
    % calculate Householder reflection of the i-th column

    for j = i+1:n
        fj = wii * R(i, j); % the head of column i is replaced by wii
        fj = fj + R(i+1:m, i)' * R(i+1:m, j); % dot of lower-part of columns i and j
        fj = fj * two_fi_inv_squared;
        % calculate columns dot

        R(i, j) = R(i, j) - fj * wii;
        R(i+1:m, j) = R(i+1:m, j) - fj * R(i+1:m, i);
        % update jth column
    end
end
% apply Householder reflections also to the other columns to the right
```

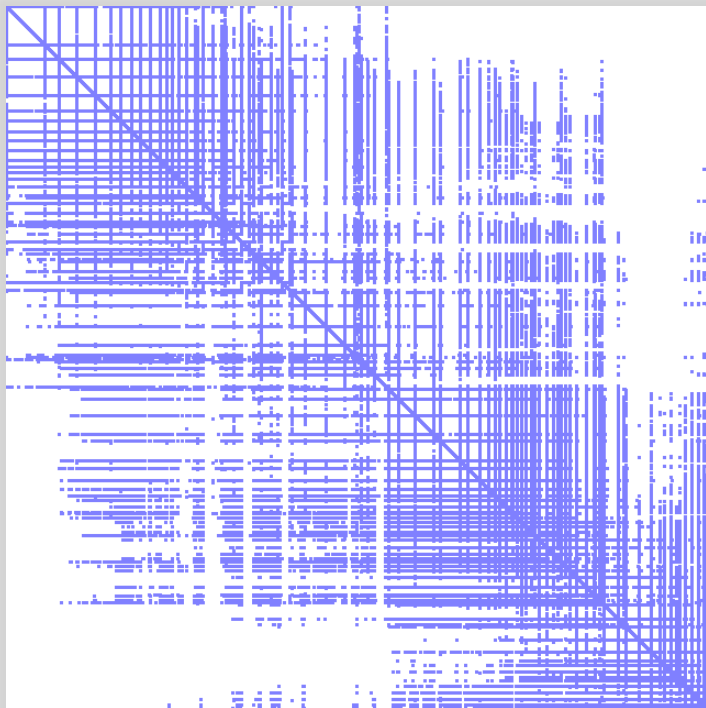
```
for j = 1:m
    fj = wii * Q(j, i); % the head of column i is replaced by wii
    fj = fj + Q(j, i+1:m) * R(i+1:m, i); % dot of lower-part of columns i and j
    fj = fj * two_fi_inv_squared;
    % calculate columns dot

    Q(j, i) = Q(j, i) - fj * wii;
    Q(j, i+1:m) = Q(j, i+1:m) - fj * R(i+1:m, i)';
    % update jth column
end
% apply Householder reflections also to the r.h.s columns, in transpose!

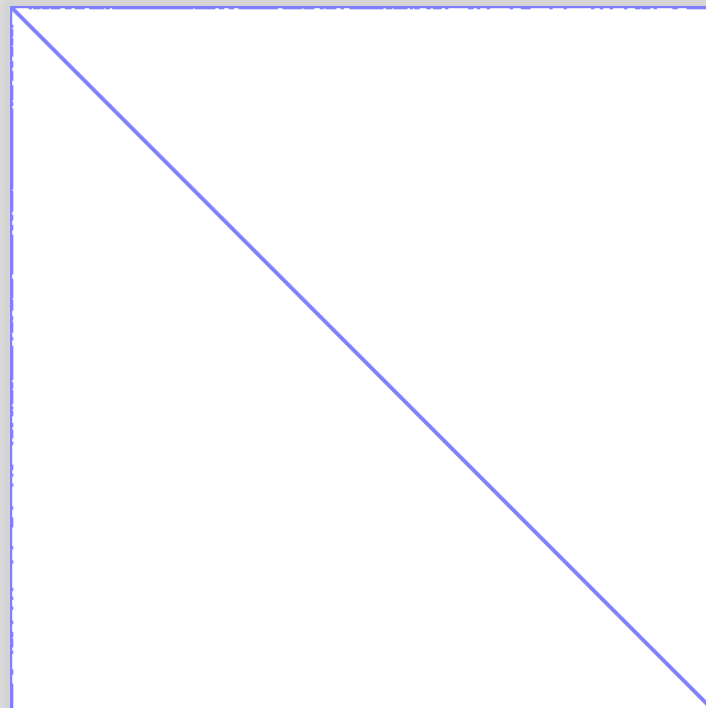
R(i, i) = dii;
R(i+1:m, i) = 0; % finally, clear the rest of the current column
end
% eliminate lower triangle of M, column by column
```

- In BA, it is possible to reorder the system to have diagonal submatrices
  - bipartite ordering, MIS

$\Lambda$



$P\Lambda P^{-1}$



92×6

57957×3

\*The sparsity is the same on the left / right. The nnz's are just very inflated to be visible and that makes the sparsity appear different.

- Then the matrix is partitioned as

$$P\Lambda P^{-1} = \begin{bmatrix} A & U \\ U^T & D \end{bmatrix}$$

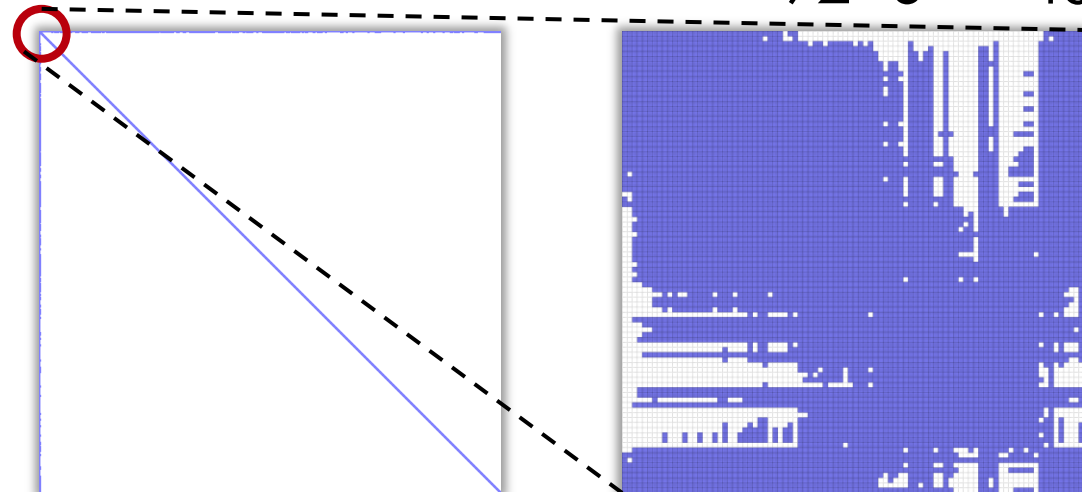
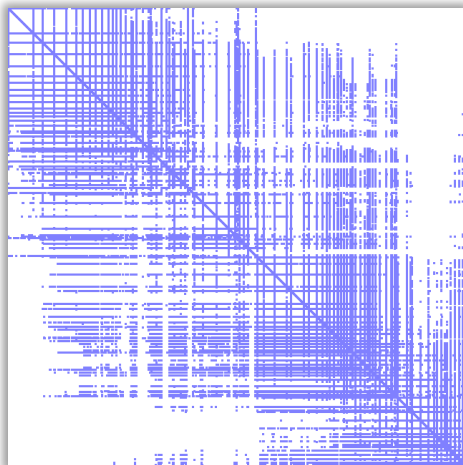
and the linear system  $\Lambda \mathbf{x} = \mathbf{b}$  partitioned as

$$\begin{bmatrix} A & U \\ U^T & D \end{bmatrix} \begin{bmatrix} \mathbf{p} \\ \mathbf{l} \end{bmatrix} = \begin{bmatrix} \mathbf{u} \\ \mathbf{v} \end{bmatrix} \text{ with } \begin{bmatrix} \mathbf{p} \\ \mathbf{l} \end{bmatrix} = P\mathbf{x}, \begin{bmatrix} \mathbf{u} \\ \mathbf{v} \end{bmatrix} = P\mathbf{b}$$

- The Schur complement of  $A$  is  $A - UD^{-1}U^T$

<1% nnz

92×6 >40% nnz



92×6

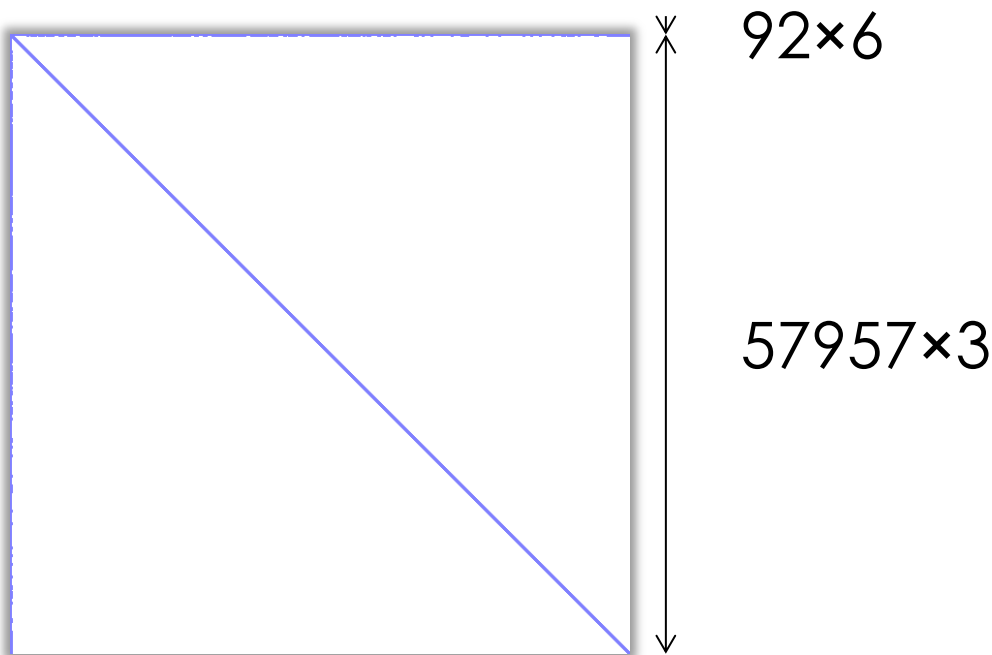
- The linear system

$$\begin{bmatrix} A & U \\ U^T & D \end{bmatrix} P^{-1} \begin{bmatrix} \mathbf{p} \\ \mathbf{l} \end{bmatrix} = P^{-1} \begin{bmatrix} \mathbf{u} \\ \mathbf{v} \end{bmatrix}$$

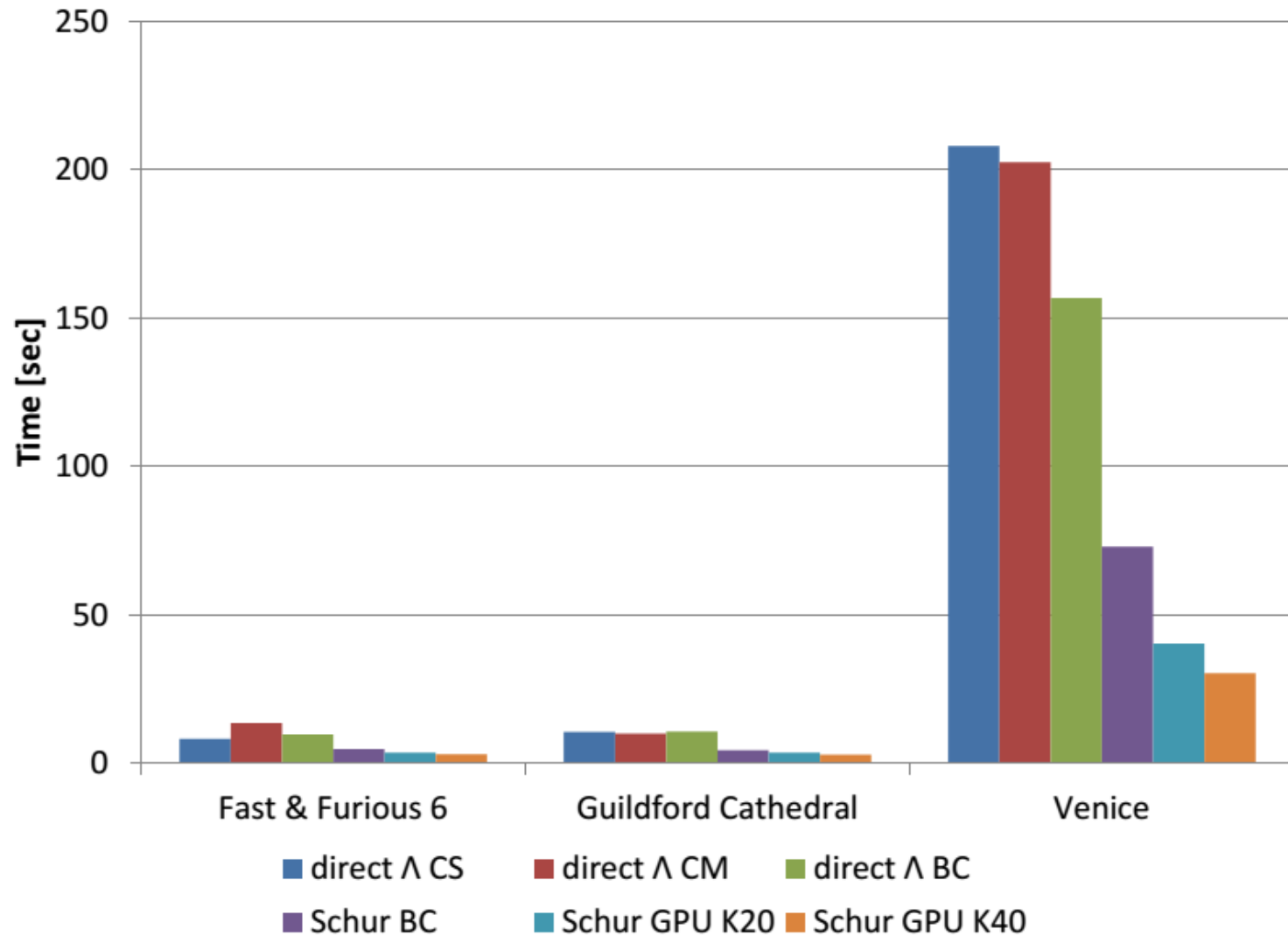
can be solved as

$$(A - UD^{-1}U^T)\mathbf{p} = \mathbf{u} - UD^{-1}\mathbf{v}$$

$$\mathbf{l} = D^{-1}(\mathbf{v} - U^T\mathbf{p})$$







Who will do seminar on

**ITERATIVE METHODS?**

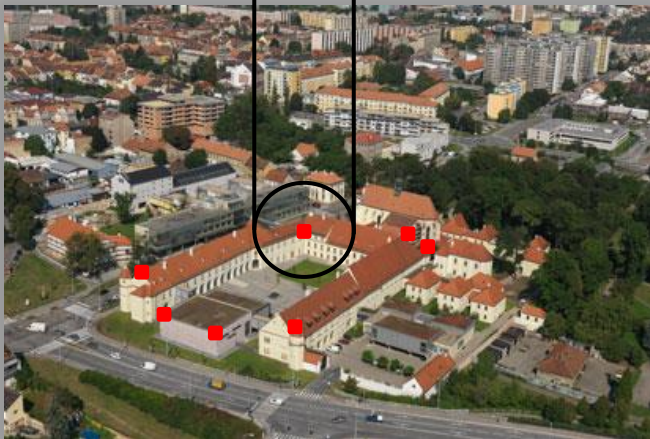
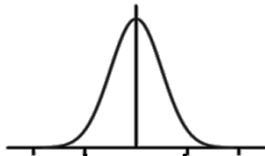
Who will do seminar on

**SUBSPACE METHODS?**

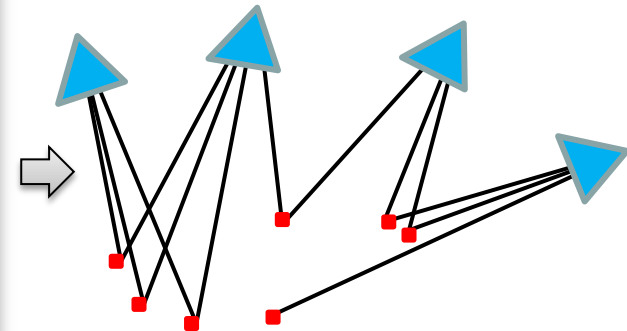
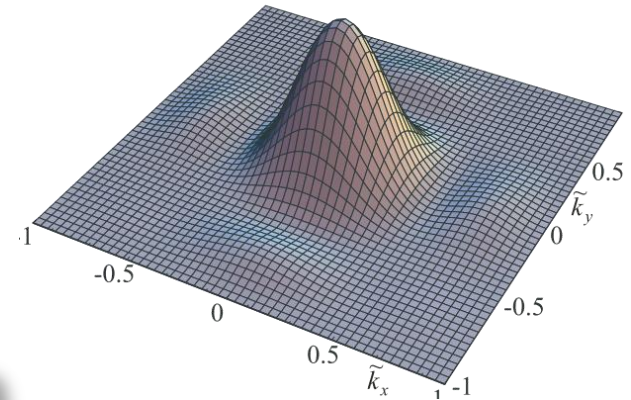
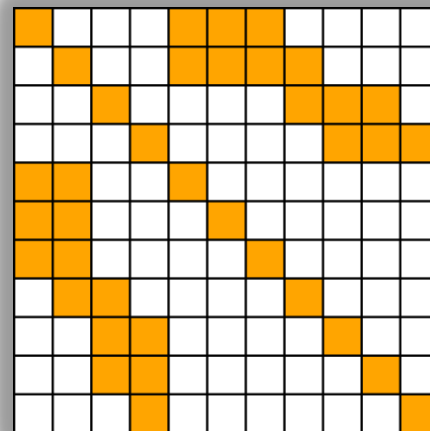
And now ...

# BACK TO BUNDLING

- System covariances (variable covariances)
  - Different from edge (observation) covariances
- Obtained by inverting the information matrix  $\Lambda$ 
  - Inverse  $\Sigma = \Lambda^{-1}$  is fully dense



$\Lambda$



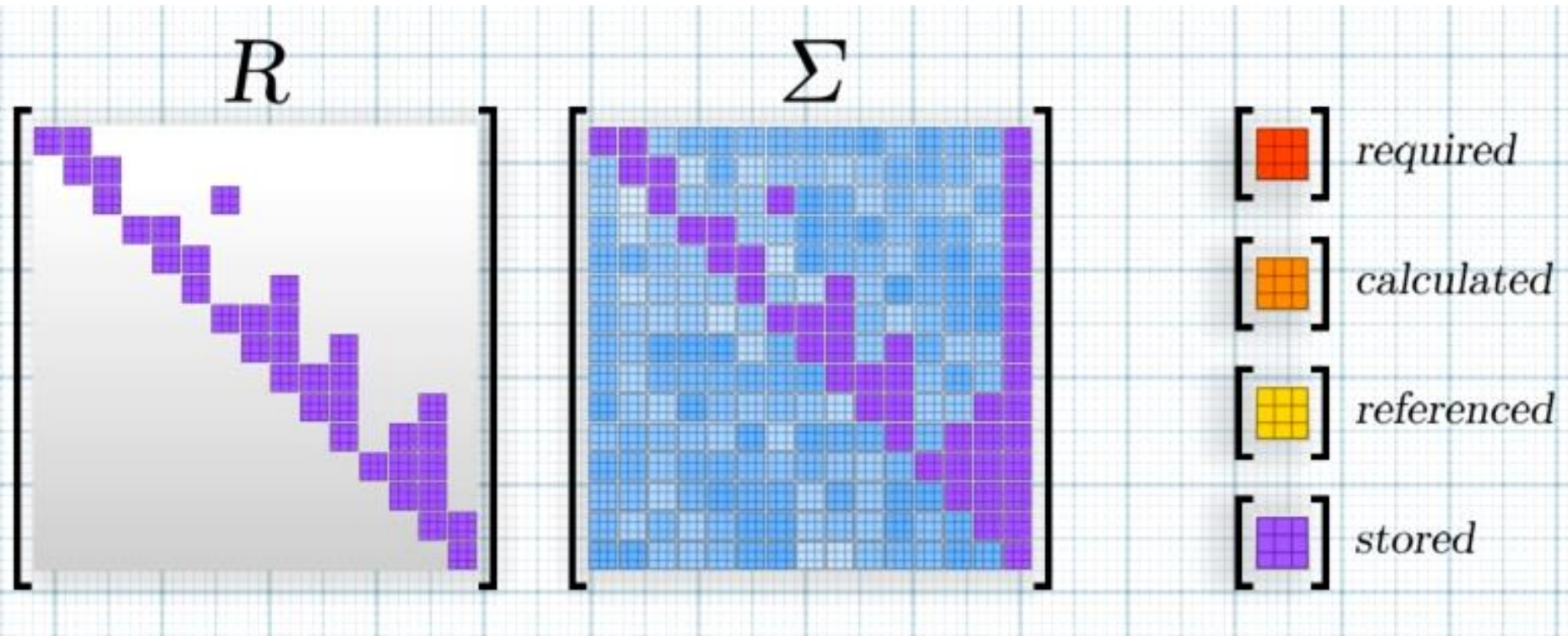
- Use Cholesky factorization & backsubstitution, keep only parts of the covariance to save memory
  - The way Google's Ceres does it, gruesome performance
- Use SVD
  - Another dead end explored by Google™

- Use recursive formula for calculating only some parts of the inverse

$$\Sigma_{ii} = \frac{1}{R_{ii}} \left( \frac{1}{R_{ii}} - \sum_{k=i+1, R_{ik} \neq 0}^n R_{ik} \Sigma_{ki} \right)$$
$$\Sigma_{ij} = \frac{1}{R_{ii}} \left( - \sum_{k=i+1, R_{ik} \neq 0}^n R_{ik} \Sigma_{kj} - \sum_{k=j+1, R_{ik} \neq 0}^n R_{ik} \Sigma_{jk} \right)$$

- [Björck, 1996, Numerical methods for least squares problems, SIAM]

- Use recursive formula for calculating only some parts of the inverse



- [Ila, 2015, Fast Covariance Recovery in Incremental Nonlinear Least Square Solvers, ICRA]

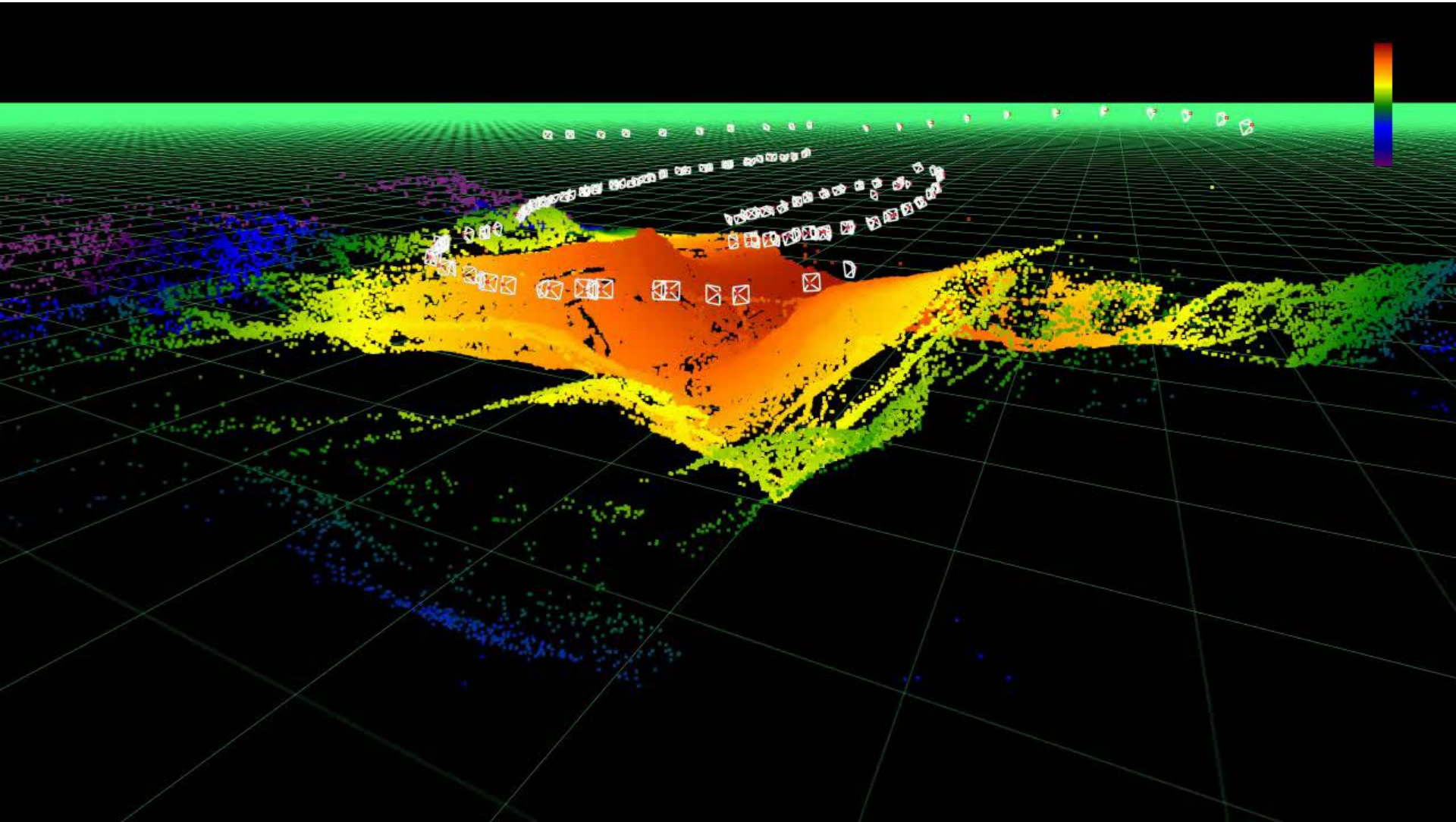


- In SLAM, we often only have small increment
- Updating covariance cheaper than recalculating from scratch

$$\begin{aligned}\Delta\Sigma &= \hat{\Sigma}A_u^T(I - A_u\hat{\Sigma}A_u^T)^{-1}A_u\hat{\Sigma} \\ \Delta\Sigma &= -\Sigma A_u^T(I + A_u\Sigma A_u^T)^{-1}A_u\Sigma\end{aligned}$$

where  $\hat{\Lambda} = \hat{A}^T \hat{A} = \Lambda + A_u^T A_u$  with  $\hat{A} = A + A_u$

- Need only a few elements of  $\hat{\Sigma}$  (can use backsubstitution)
- [Ila, 2015, Fast Covariance Recovery in Incremental Nonlinear Least Square Solvers, ICRA]





The end

# ARE YOU STILL ALIVE?