# Numerical Methods I
## MATH-GA 2010.001/CSCI-GA 2420.001

Benjamin Peherstorfer
Courant Institute, NYU

Based on slides by G. Stadler and A. Donev

# Today

**Last time**

- ▶ Solving systems of linear equations
- ▶ LU decomposition

**Today**

- ▶ Recap: linear systems. Sparse matrices
- ▶ Linear least square problems
- ▶ Geometric perspective on the normal equations
- ▶ Orthogonalization with Gram-Schmidt

**Announcements**

- ▶ Homework 2 has been posted, due Tue (!), Oct 11 before class
- ▶ Final exam on Mon, Dec 19, 5.10 - 7.00pm
- • Office hour : Tue, 1pm

# Gauss elimination and LU factorization

$$\begin{bmatrix} a_{11}^{(1)} & a_{12}^{(1)} & a_{13}^{(1)} \\ a_{21}^{(1)} & a_{21}^{(1)} & a_{23}^{(1)} \\ a_{31}^{(1)} & a_{32}^{(1)} & a_{33}^{(1)} \end{bmatrix} \begin{bmatrix} x_1 \\ x_2 \\ x_3 \end{bmatrix} = \begin{bmatrix} b_1^{(1)} \\ b_2^{(1)} \\ b_3^{(1)} \end{bmatrix}$$

*index of step*

multiply first row by $l_{i1}$

and subtract from 2nd row

$$l_{21} = \frac{a_{21}^{(1)}}{a_{11}^{(1)}}$$

$$\begin{bmatrix} a_{11}^{(1)} & a_{12}^{(1)} & a_{13}^{(1)} \\ 0 & a_{22}^{(2)} & a_{23}^{(1)} \\ 0 & a_{32}^{(2)} & a_{33}^{(2)} \end{bmatrix} \begin{bmatrix} x_1 \\ x_2 \\ x_3 \end{bmatrix} = \begin{bmatrix} b_1^{(1)} \\ b_2^{(2)} \\ b_3^{(2)} \end{bmatrix}$$

eliminate using

$$l_{32} = \frac{a_{32}^{(2)}}{a_{22}^{(2)}}$$

$$\begin{bmatrix} a_{11}^{(1)} & a_{12}^{(1)} & a_{12}^{(1)} \\ 0 & a_{22}^{(2)} & a_{23}^{(2)} \\ 0 & 0 & a_{33}^{(3)} \end{bmatrix} \begin{bmatrix} x_1 \\ x_2 \\ x_3 \end{bmatrix} = \begin{bmatrix} b_1^{(1)} \\ b_2^{(2)} \\ b_3^{(3)} \end{bmatrix}$$

Idea: Store multipliers $l_{21}, l_{31}, l_{32}$ in matrix:

$$\begin{bmatrix} u_{11} & u_{12} & u_{13} \\ l_{21} & u_{22} & u_{23} \\ l_{31} & l_{32} & u_{33} \end{bmatrix}$$

$\begin{bmatrix} L & U \end{bmatrix}$

$$A = LU, \quad L = \begin{bmatrix} 1 & & \\ l_{21} & 1 & \\ l_{31} & l_{31} & 1 \end{bmatrix}, \quad U = \begin{bmatrix} u_{11} & \cdots & u_{13} \\ 0 & & \\ 0 & 0 & \ddots \end{bmatrix}$$

# LU with pivoting

Zero diagonal entries (pivots) pose a problem $\rightsquigarrow$ pivoting by swapping rows

$\rightsquigarrow$ board
$$\begin{bmatrix} 1 & 1 & 3 \\ 2 & 2 & 2 \\ 3 & 6 & 4 \end{bmatrix}$$

$\xrightarrow[\text{elimination}]{\text{Gauss}}$
$$\begin{bmatrix} 1 & 1 & 3 \\ l_{21}=2 & 0 & -4 \\ l_{31}=3 & 3 & -5 \end{bmatrix}$$

Problem!
Cannot divide
by zero,
pivot element is
zero.

Swap rows
$\xrightarrow{\ \ }$
2 & 3
$$\begin{bmatrix} 1 & 1 & 3 \\ 3 & 3 & -5 \\ 2 & 0 & -4 \end{bmatrix}$$

done — In general, use 3 as pivot.

$$\boxed{LU = PA}$$

Theorem: This always exists with
appropriate perm. matrix P

P.... permutation matrix.

$P_2 \begin{bmatrix} 1 & 0 & 0 \\ 0 & 0 & 1 \\ 0 & 1 & 0 \end{bmatrix}$

# Recap

▶ For any square (regular or singular) matrix $A$, partial (row) pivoting ensures exists of

$$PA = LU$$

where $P$ is a permutation matrix

▶ Pivoting (w.r.t. $\max |a_{ij}|$) leads to a backward stable algorithm.

▶ There also is full pivoting (rows + columns)

$$PAQ = LU$$

to further increases stability but it usually is not worth it in practice (higher costs to search for pivoting element over rows and columns but little improvement in terms of stability)

# Recap: Solving linear systems

▶ Once an LU factorization is available, solving a linear system is cheap:

$$\underline{Ax} = \underline{LU}x = \underline{L}(\underline{Ux}) = \underline{Ly} = \underline{b}$$

▶ Solve for $y$ using *forward substitution*
▶ Solve for $x$ by using *backward substitution* $Ux = y$

Forward/backward substitution requires

$$\frac{n(n-1)}{2} \text{ multiplications/additions,}$$

$$n \text{ divisons.}$$

Overall: $\sim n^2$ floating point operations (flops) $\rightsquigarrow$ costs scale as $\mathcal{O}(n^2)$

# Recap: Costs

For forward [backward] substitution at step $k$ there are $\approx k \ [(n-k)]$ multiplications and subtractions plus a few divisions. The total over all $n$ steps is

$$\sum_{k=1}^{n} k \in \mathcal{O}(n^2)$$

$\rightsquigarrow$ the number of floating-point operations (FLOPs) scales as $\mathcal{O}(n^2)$

For Gaussian elimination, at step $k$, there are $\approx (n-k)^2$ operations. Thus, the total scales as

$$\sum_{k=1}^{n}(n-k)^2 \in \mathcal{O}(n^3)$$

Summary:
- ▶ Directly applying Gaussian elimination (=LU + fwd/bwd) scales as $\mathcal{O}(n^3)$
- ▶ Computing LU decomposition scales as $\mathcal{O}(n^3)$
- ▶ Forward/backward substitution scales as $\mathcal{O}(n^2)$
- ▶ LU + forward/backward scales as $\mathcal{O}(n^3)$ $\rightsquigarrow$ can *reuse LU* for other $b$

# Recap: Summary

▶ The condition of solving a linear system $Ax = b$ is determined by the condition number of the matrix $A$

$$\kappa(A) = \|A\|\|A^{-1}\| \geq 1$$

▶ Gaussian elimination can be used to solve general square linear systems and produces a factorization, if it exists

$$A = LU$$

▶ Partial pivoting is sufficient for existence and stability of the LU decomposition

$$PA = LU, \qquad A = \tilde{L}U$$

▶ The Cholesky factorization $A = LL^T$ exists if $A$ is spd and then it is the better choice (stabler, cheaper) than LU

▶ Rely on the highly optimized routines in Matlab (LAPACK) and other software packages than implementing these algorithms yourself ⇝ take the course on HPC next spring to learn more about the efficient *implementation* of these algorithms

# Sparse matrix

▶ A matrix where a substantial fraction of the entries are zero is called a sparse matrix. Typically, only $\mathcal{O}(N)$ non-zero entries are allowed.

▶ If we have only $\mathcal{O}(N)$ non-zero entries, then store only those; in contrast to dense matrices. Exploiting sparsity is important (life saving) for large matrices

▶ The structure of a sparse matrix refers to the set of indices $i, j$ such that $|a_{ij}| > 0$ and is visualized in Matlab with spy

▶ The sparsity structure is the most important property that determines whether an efficient method exists

▶ The structure of sparse matrices comes from the nature of the problem

# Banded matrices

Banded matrices are a very special but common type of sparse matrices, e.g., tridiagonal matrices

$$\begin{bmatrix} a_1 & c_1 & & & 0 \\ b_2 & a_2 & & \ddots & \\ & \ddots & \ddots & & c_{n-1} \\ 0 & & & b_n & a_n \end{bmatrix}$$

For example, think of the Laplace problem $u''(x) = f(x)$ on the unit interval and a finite-difference discretization

$$u''(x) \approx \frac{u(x+h) - 2u(x) + u(x-h)}{h^2}$$

on an equidistant grid. This leads to a system of equations with a tridiagonal matrix $\rightsquigarrow$ Numerical Methods II (Spring semester)

There exist special techniques for banded matrices that are much faster than the general case, e.g., only $8n$ FLOPs

# Decomposing sparse matrices

There also are general methods for dealing with sparse matrices, such as the sparse LU factorization.

How well they work depends on the structure of the matrix. What could go wrong?

# Decomposing sparse matrices

There also are general methods for dealing with sparse matrices, such as the sparse LU factorization.

How well they work depends on the structure of the matrix. What could go wrong?

When factorizing sparse matrices, the factors, e.g., $L$ and $U$, can be much less sparse than $A \rightsquigarrow$ fill-in

For many sparse matrices, there is a large fill-in

► Pivoting can help to reduce fill-in
► However, often "good" pivoting for sparsity leads to less stable behavior and vice versa

# Sparse matrices in Matlab

```
 1: % S = sparse(i,j,v) generates a sparse matrix S
 2: % from the triplets i, j, v with S(i(k),j(k)) = v(k).
 3: >> A = sparse([1 2 2 4 4], [3 1 4 2 3], 1:5)
 4: A =
 5:     (2,1)        2
 6:     (4,2)        4
 7:     (1,3)        1
 8:     (4,3)        5
 9:     (2,4)        3
10: >> whos A
11:    Name          Size                  Bytes   Class
          Attributes
12:    A             4x4                     120   double    sparse
13: >> nnz(A)
14: ans =
15:        5
```
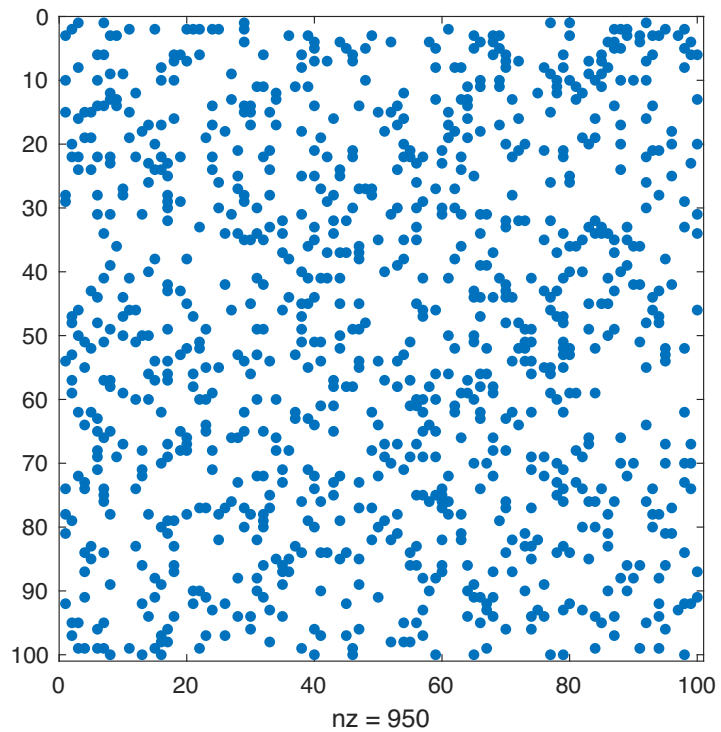
```
 1: % S = sparse(i,j,v,m,n,nz) allocates space for nz
      nonzero elements.
 2: % Use this syntax to allocate extra space for nonzero
      values to be filled in after construction.
 3: >> A = sparse([], [], [], 4, 4, 5);
 4: >> A(2, 1) = 2; A(4, 2) = 4; A(1, 3) = 1; A(4, 3) = 5;
      A(2, 4) = 3;
 5: >> full(A)
 6:
 7: ans =
 8:
 9:        0        0        1        0
10:        2        0        0        3
11:        0        0        0        0
12:        0        4        5        0
```
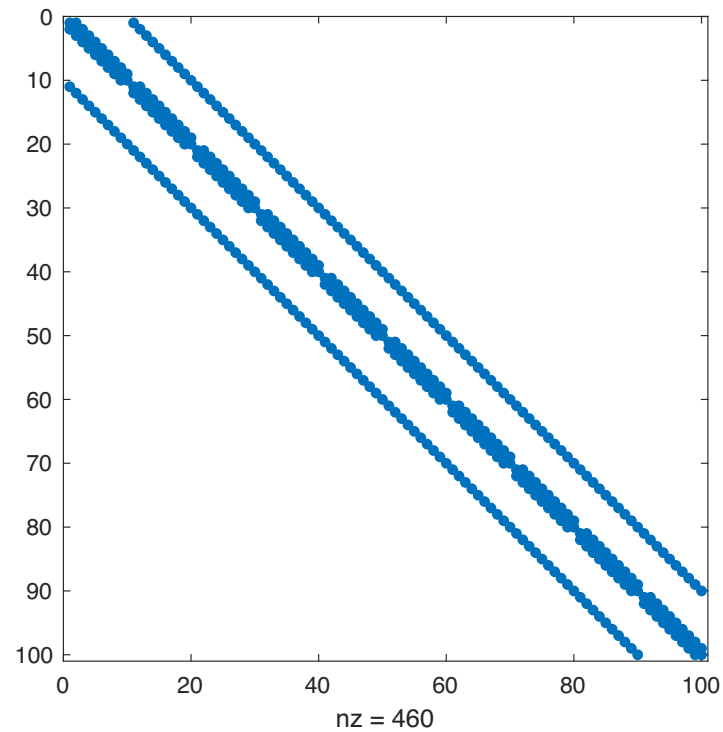
```matlab
1: % generate a random sparse matrix with density 10% and
       size 100x100
2: >> B = sprand(100, 100, 0.1);
3: % the sparse block tridiagonal matrix of order n^2
       resulting from discretizing Poisson's equation with
       the 5-point operator on an n-by-n mesh.
4: >> X = gallery('poisson', 10);
5: >> spy(B);
6: >> spy(X);
```

(a) matrix $B$
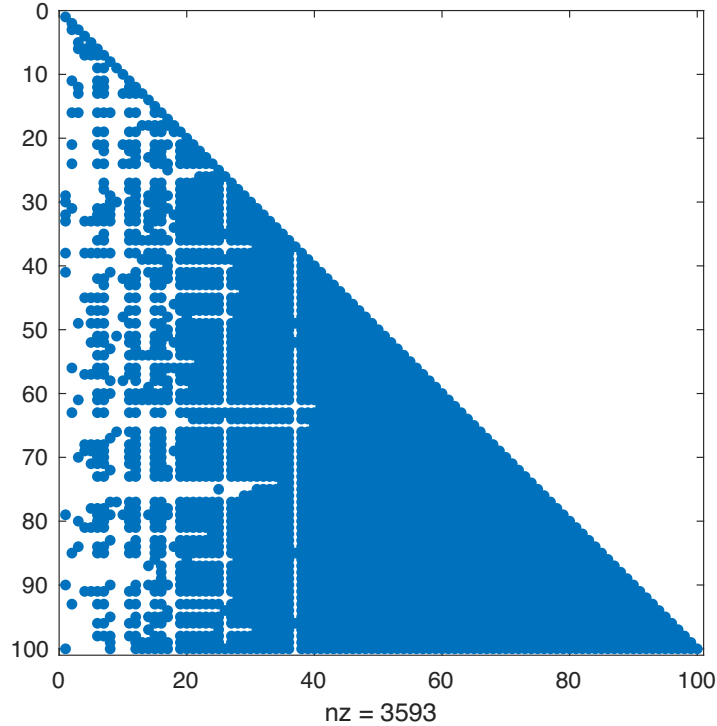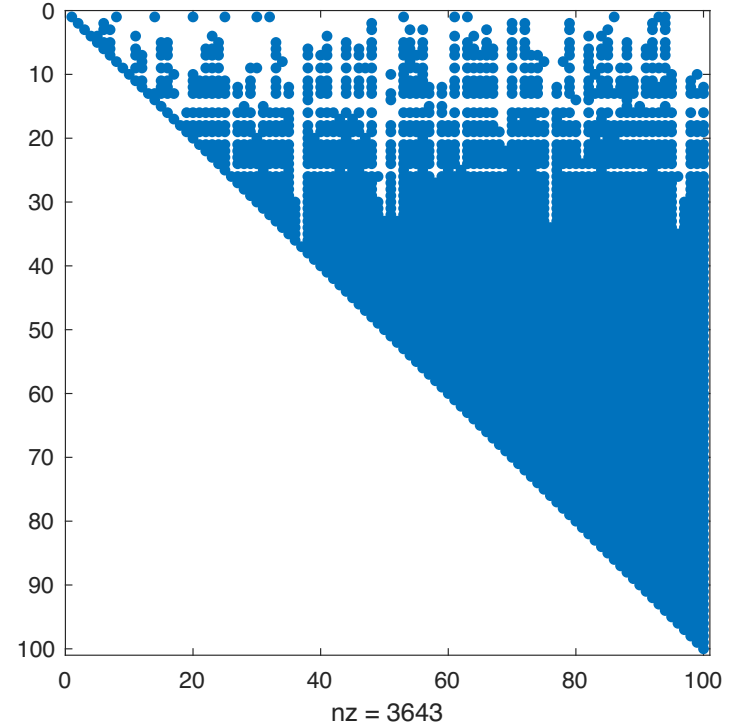
(b) matrix $X$

```
1: >> [L, U, P] = lu(B);
2: >> spy(L);
3: >> spy(U);
4:
5: >> [L, U, P] = lu(X);
6: >> spy(L);
7: >> spy(U);
```
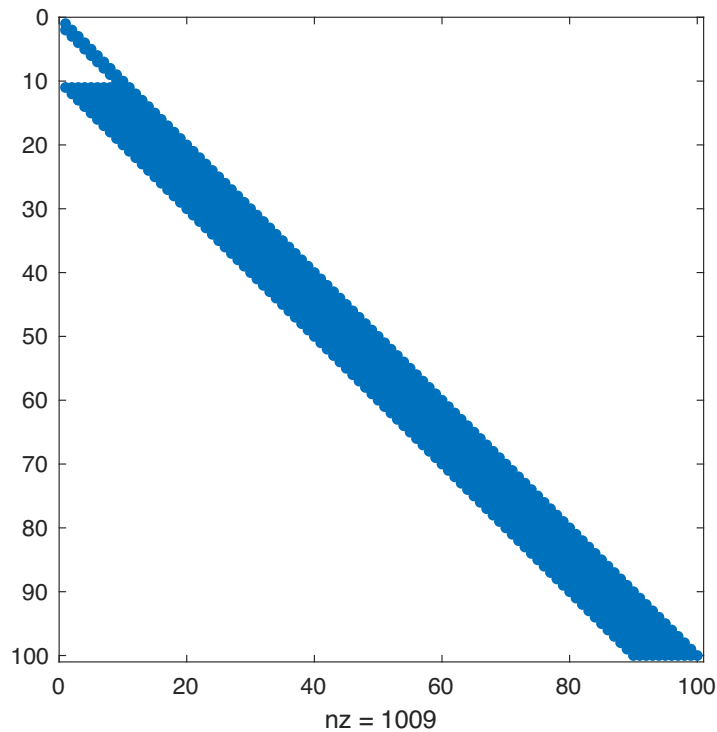
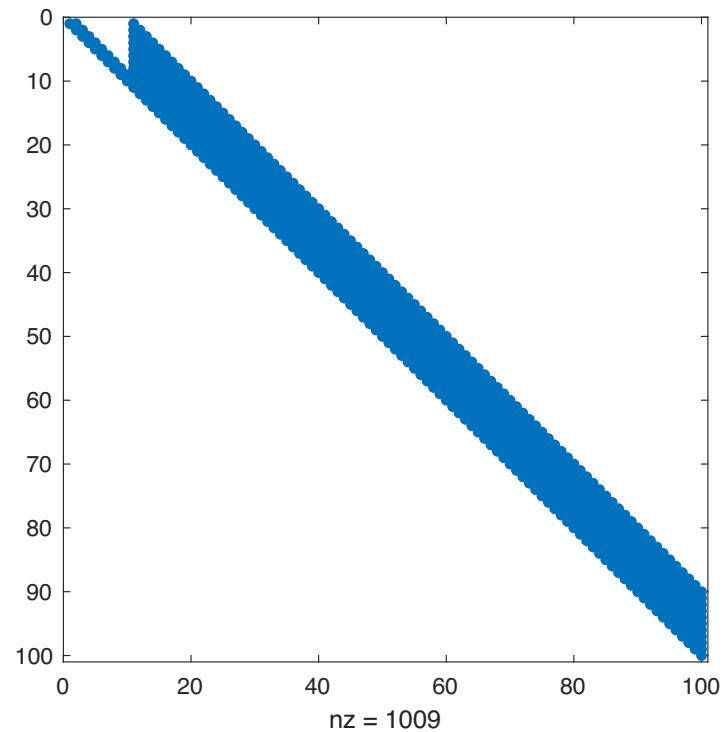(a) factor $L$ of $B$          (b) factor $U$ of $B$

A lot of fill-in! Factors $L$ and $U$ are not sparse, even though matrix $B$ is sparse

(a) factor $L$ of $X$           (b) factor $U$ of $X$

Though better than for matrix $B$ with random sparsity structure, there still are many more non-zero entries in the factors of $X$ than in $X$ itself.
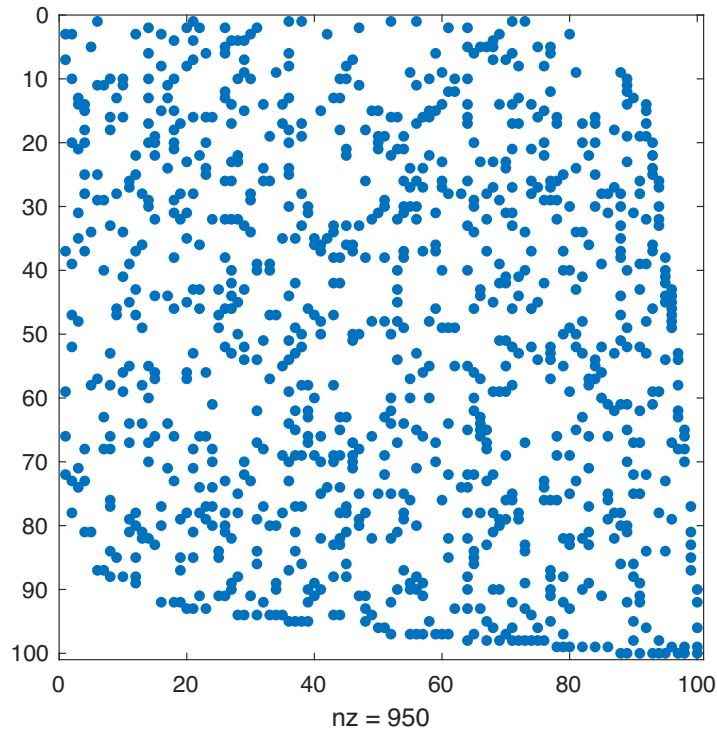
Changing the sparsity structure via re-ordering the matrix can help to reduce fill-in. For example, in Matlab the sparse reverse Cuthill-McKee ordering is implemented.

The re-ordered matrix tends to have its nonzero elements closer to the diagonal. This is a good preordering for LU or Cholesky factorization of matrices.
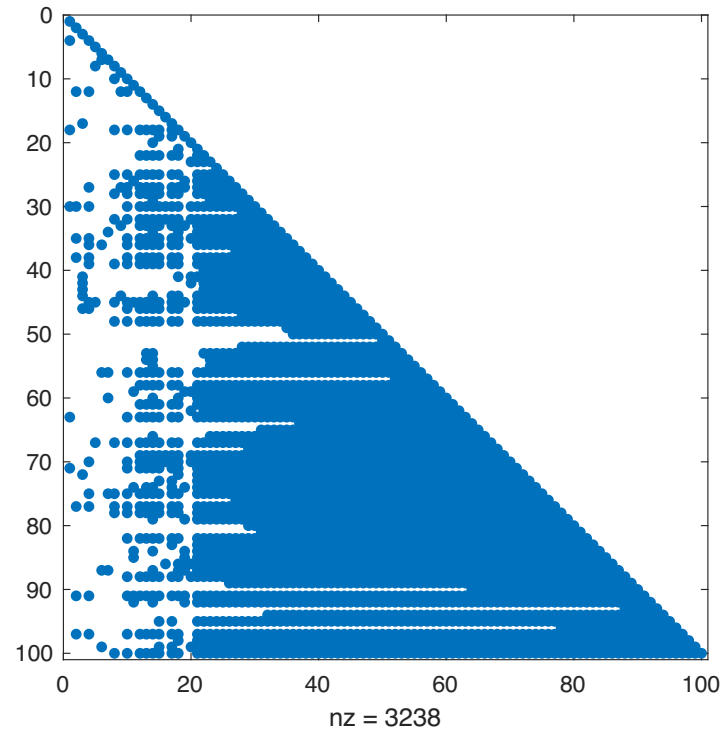
```
 1: >> p = symrcm(B);
 2: >> spy(B(p, p));
 3: >> [L, U, P] = lu(B(p, p));
 4: >> spy(L);
 5: >> spy(U);
 6:
 7: >> p = symrcm(X)
 8: >> spy(X(p, p));
 9: >> [L, U, P] = lu(X(p, p));
10: >> spy(L);
11: >> spy(U);
```
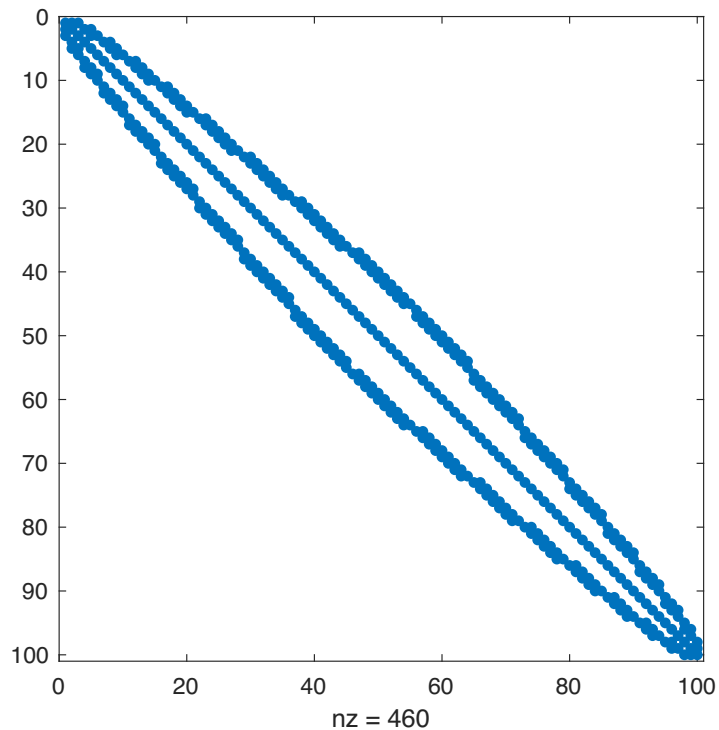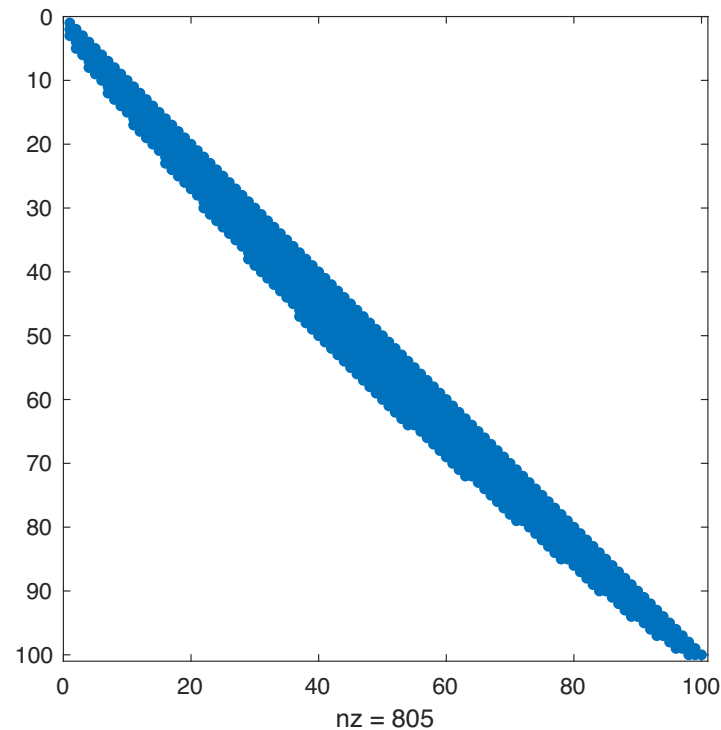
(a) re-ordered $B$ matrix　　　　　　　(b) factor $L$ of re-ordered $B$

▶ Notice how the non-zero elements tend to be closer to the diagonal of the re-ordered $B$ compared to the original $B$

▶ The fill-in is reduced from $\sim 3500$ to $\sim 3200$ non-zero entries; reduction of $< 10\%$

▶ It is hard to find a good ordering for matrix with a random sparsity structure

(a) re-ordered $X$ matrix



(b) factor $L$ of re-ordered $X$

► Non-zero entries of re-ordered $X$ are closer to the diagonal than for the original $X$

► Fill-in is reduced by a roughly 20% from $\sim 1000$ to $\sim 800$ non-zero entries

# Conclusions/summary

While there are general techniques for dealing with sparse matrices that help greatly, it all depends on the structure of the matrix

Pivoting has a dual, sometimes conflicting goal:

1. Reduce fill-in, i.e., improve memory use: Still active subject of research!
2. Reduce roundoff errors, i.e., improve stability. Typically some threshold pivoting is used only when needed

For many sparse matrices *iterative methods* (later) are required when large fill-in.

# Least-squares problems

# Least-squares problems

Given data points/measurements

$$(t_i, b_i), \quad i = 1, \ldots, m$$

and a model function $\phi$ that relates $t$ and $b$:

$$b = \phi(t; x_1, \ldots, x_n),$$

where $x_1, \ldots, x_n$ are model function parameters. If the model is supposed to describe the data, the deviations/errors

$$\Delta_i = b_i - \phi(t_i, x_1, \ldots, x_n)$$

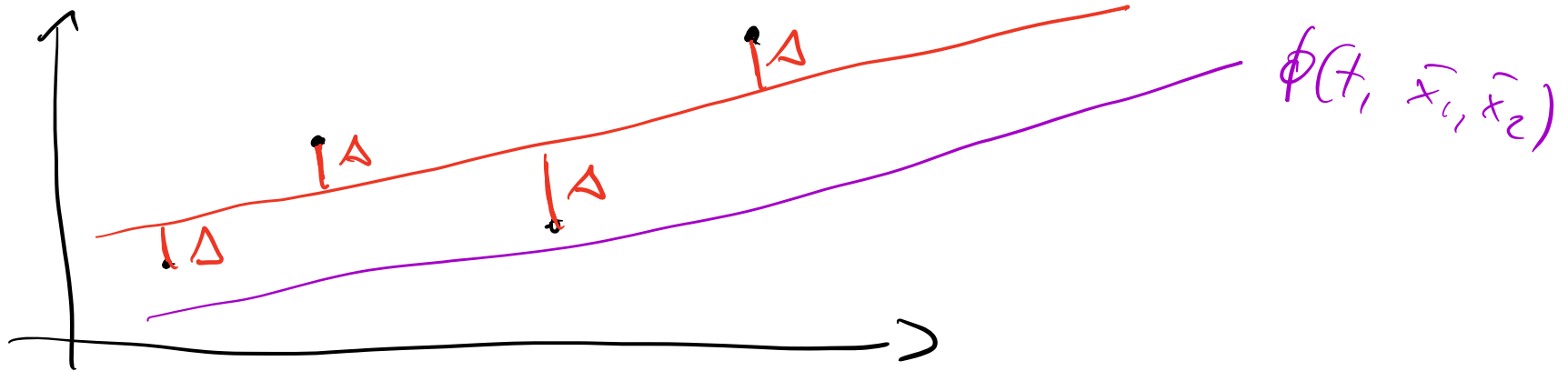should be small. Thus, to fit the model to the measurements, one must choose $x_1, \ldots, x_n$ appropriately.

# Least-squares problems

$$\text{data} \quad \{(t_i, b_i)\}_{i=1}^{m}$$

$$\text{parametrization}$$

$$b = \phi(t; x_1, \ldots, x_n)$$

$$\Delta_i = b_i - \phi(t_i, x_1, \ldots, x_n)$$



$\phi(t, \tilde{x}_1, \tilde{x}_2)$

# Least-squares problems

Least squares: Find $x_1, \ldots, x_n$ such that

$$\frac{1}{2} \sum_{i=1}^{m} \Delta_i^2 \to \min$$

Weighted least squares: Find $x_1, \ldots, x_n$ such that

$$\frac{1}{2} \sum_{i=1}^{m} \left( \frac{\Delta_i}{\delta b_i} \right)^2 \to \min,$$

where $\delta b_i > 0$ contain information about how much we trust the $i$th data point.

# Least-squares problems (cont'd)

Alternatives to using squares:

$L^1$ error: Find $x_1, \ldots, x_n$ such that

$$\sum_{i=1}^{m} |\Delta_i| \to \min$$

Result can be very different, other statistical interpretation, more stable with respect to outliers.

$L^\infty$ error: Find $x_1, \ldots, x_n$ such that

$$\max_{1 \le i \le m} |\Delta_i| \to \min$$

Keeps the worst-case error small (risk averse)

# Linear least-squares

We assume (for now) that the model depends linearly on $x_1, \ldots, x_n$, e.g.:

$$\phi(t; x_1, \ldots x_n) = a_1(t)x_1 + \ldots + a_n(t)x_n$$

$\rightsquigarrow$ board

$$\phi(t; x_1, x_2) = x_1 t^2 + x_2 \exp(t)$$

$$\Delta_1 = b_1 - (x_1 t_1^2 + x_2 \exp(t_1))$$

$$\Delta_2 = b_2 - (x_1 t_2^2 + x_2 \exp(t_2))$$

$$\Delta_3 = b_3 - (x_1 t_3^2 + x_2 \exp(t_3))$$

$$\min_{x_1, x_2} \sum_{i=1}^{3} \Delta_i^2 \rightsquigarrow \text{linear lsl sq problem}$$

# Linear least-squares

Choosing the least square error, this results in

$$\min_{\boldsymbol{x}} \|A\boldsymbol{x} - \boldsymbol{b}\|^2,$$

where $\boldsymbol{x} = (x_1, \ldots, x_n)^T$, $\boldsymbol{b} = (b_1, \ldots, b_m)^T$, and $a_{ij} = a_j(t_i)$.

In the following, we study the overdetermined case, i.e., $m \geq n \rightsquigarrow$ board

$$A = \begin{bmatrix} t_1^2 & \exp(t_1) \\ t_2^2 & \exp(t_2) \\ t_3^2 & \exp(t_3) \end{bmatrix} / \quad b = \begin{bmatrix} b_1 \\ b_2 \\ b_3 \end{bmatrix} / \quad x = \begin{bmatrix} x_1 \\ x_2 \end{bmatrix}$$

$$\in \mathbb{R}^{3 \times 2} \qquad\qquad \in \mathbb{R}^3 \qquad\qquad \in \mathbb{R}^2$$

# Linear least-squares

Different perspective:

Consider non-square matrices $A \in \mathbb{R}^{m \times n}$ with $m \geq n$ and $\text{rank}(A) = n$. Then the system

$$Ax = b$$

does, in general, not have a solution (more equations than unknowns). We thus instead solve a minimization problem

$$\min_x \|Ax - b\|^2 = \min_x \Phi(x)$$

How can we solve this optimization problem?

Because we consider the Euclidean norm $\|.\|_2$, we obtain

$$\Phi(x) = (Ax - b)^T (Ax - b) = x^T A^T Ax - \ldots$$

which is quadratic in $x$ if $A$ has full rank $\leadsto$ convex in $x$

Therefore, the critical point is the global optimum

$$\nabla \Phi(x) = A^T (2(Ax - b)) = 0$$

which satisfies the *normal equations*

$$A^T Ax = A^T b.$$

If $A$ is full rank, $\text{rank}(A) = n$, then $A^T A$ is positive definite and the normal equations can be solved with the Cholesky factorization (warning: we shouldn't do this, can you already see why?)
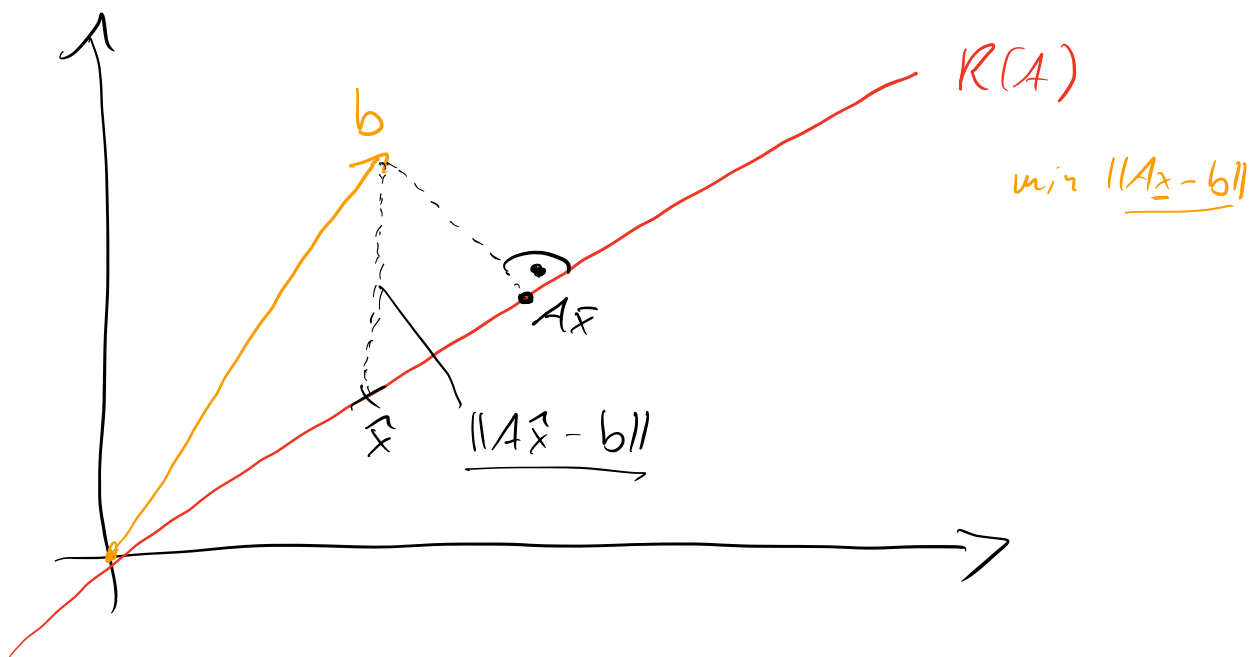
# A geometry perspective on the normal equations

$$\min_{x \in \mathbb{R}^m} \| Ax - b \|_2^2$$

Geometric sketch

$$\bar{x} = \arg\min_{x \in \mathbb{R}^m} \| Ax - b \|_2^2$$

$$\Downarrow$$

$$A\bar{x} - b \perp R(A) \longleftarrow \{ Ax \mid x \in \mathbb{R}^m \}$$

$R(A)$

$\min \|A\underline{x} - b\|$

$b$

$A\hat{x}$

$\hat{x}$

$\|A\hat{x} - b\|$

Orthogonal projection

Finite-dim space $\underline{V}$ with $\langle \cdot, \cdot \rangle$ and subspace $\underline{U} \subset V$. Let $U^\perp$ be the orthogonal complement w.r.t. $V$

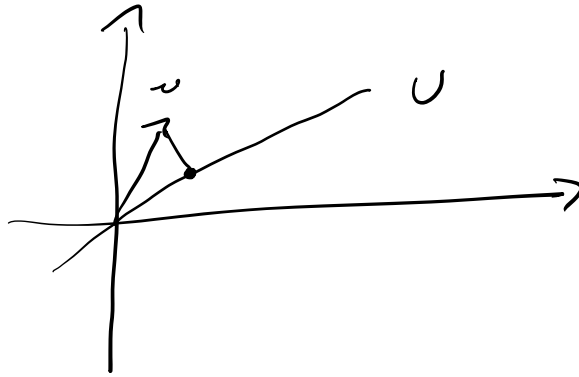$$U^\perp = \{\underline{v} \in V \mid \langle v, u \rangle = 0 \quad \forall u \in U\}$$

Then

$$\arg\min_{u' \in U} \|u' - \underline{v}\| = u \iff v - u \in U^\perp$$

Proof: let $u \in U$ such that $\underline{v - u \in U^\perp}$
for all $u' \in U$

$$\|v - u'\|^2 = \|v - u + u - u'\|^2$$
$$= \|v - u\|^2 + \|u - u'\|^2 \geq \|v - u\|^2$$
$$\text{iff } u \neq u'$$

$$u = P(v)$$

P is linear in v

thus

$$u = Pv$$

$$\bar{x} = \underset{x \in \mathbb{R}^n}{\arg\min} \; \|Ax - b\|_2^2$$

$$V = \mathbb{R}^m, \quad U = R(A) \subset V$$

$$\|b - Ax\|^2 = \min \iff \langle b - Ax, Ax'\rangle = 0 \quad \forall x' \in \mathbb{R}^n$$

$$\iff \langle A^T(b - Ax), \underline{x}'\rangle = 0$$

$$\iff \boxed{A^T b = A^T Ax}$$

$$\iff \underline{x} = (A^T A)^{-1} A^T b$$

$$Ax = Pb = A(A^T A)^{-1} A^T b \qquad x = (A^T A)^{-1} A^T b$$

# Linear least-squares problems

Solving the normal equations

$$A^T A \bar{\mathbf{x}} = A^T \mathbf{b}$$

requires:

- ▶ computing $A^T A$ (which is $O(mn^2)$)
- ▶ condition number of $\underline{A^T A}$ is square of condition number of $A$; (problematic for the Choleski factorization)

⇝ derive condition of $A^T A$ on board

$$k_2(A^\top A) = k_2(A)^2$$

$$k_2(A)^2 = \|A\|_2^2 \, \|A^{-1}\|_2^2 = \frac{\max\limits_{\|x\|=1} \|Ax\|_2^2}{\min\limits_{\|x\|=1} \|Ax\|_2^2}$$

$$= \frac{\max\limits_{\|x\|=1} \langle Ax, Ax \rangle}{\min\limits_{\|x\|=1} \langle Ax, Ax \rangle}$$

$$= \frac{\max\limits_{\|x\|=1} \boxed{\langle A^\top A x, x \rangle}}{\min\limits_{\|x\|=1} \langle A^\top A x, x \rangle}$$

$$= \frac{\lambda_{max}(A^\top A)}{\lambda_{min}(A^\top A)}$$

$A^\top A$ is spd

$$\lambda_{max}(A^\top A) = \sqrt{\lambda_{max}(A^\top A)^2} \overset{spd}{=} \sqrt{\lambda_{max}((A^\top A)^2)}$$

$$= \sqrt{\lambda_{max}((A^\top A)^\top (A^\top A))}$$

$$= \|A^\top A\|_2$$

$$\underline{\quad 1 \quad}_{\textstyle \lambda_{min}(A^\top A)} = \|(A^\top A)^{-1}\|_2$$

$$\Rightarrow k_2^2(A) = \|A^\top A\|_2 \, \|(A^\top A)^{-1}\|_2 = k_2(A^\top A)$$

# Linear least-squares problems

Conditioning

Solving the normal equation is equivalent to computing $P\boldsymbol{b}$, the orthogonal projection of $\boldsymbol{b}$ onto the subspace $V$ spanned by columns of $A$.

Let $P : \mathbb{R}^m \to V$ be an orthogonal projection onto $V \subseteq \mathbb{R}^n$. For $b \in \mathbb{R}^m$, denote by $\theta$ the angle between $b$ and $V$ defined by

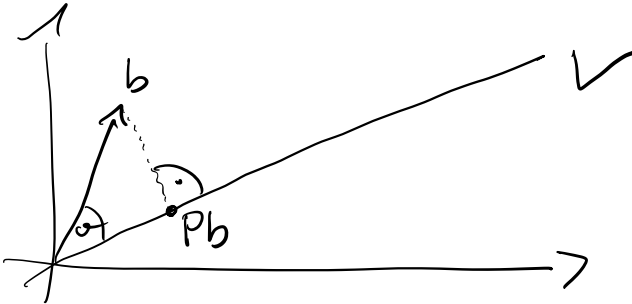$$\sin(\theta) = \frac{\|b - Pb\|_2}{\|b\|_2} .$$

The the relative condition number of projecting $b$ onto $V$ with $P$ with respect to the 2-norm is
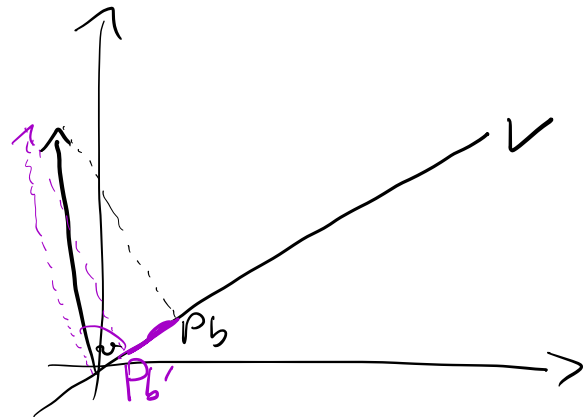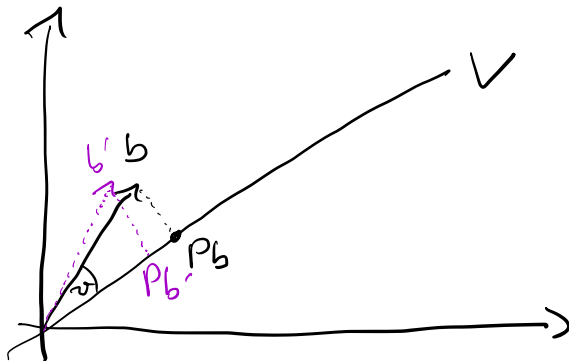
$$\kappa_{\text{rel}}(b) = \frac{1}{\cos(\theta(b))} \|P\|_2 .$$

$$P : \mathbb{R}^m \to V$$

$$\sin(v) = \frac{\|b - Pb\|_2}{\|b\|_2}$$



Relative condition number of $(P, b)$ w.r.t. 2-norm

$$k_{rel} = \boxed{\frac{1}{\cos(v)}} \|P\|_2$$

$$k_{rel} = \frac{||x||}{||f(x)||} \, || f'(x)||$$

$$k_{rel} = \frac{||b||}{||Pb||} \, || P'(b)||$$

(1) $P'(b) = P$

(2) $\dfrac{|| Pb||^2}{|| b||^2} =$

$Pb \perp b - Pb$

$$||b||^2 = || Pb + b - Pb||^2$$
$$= ||Pb||^2 + ||b - Pb||^2$$
$$||Pb||^2 = ||b||^2 - ||b - Pb||^2$$
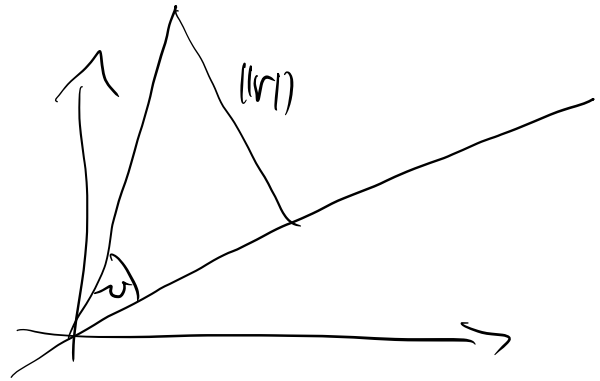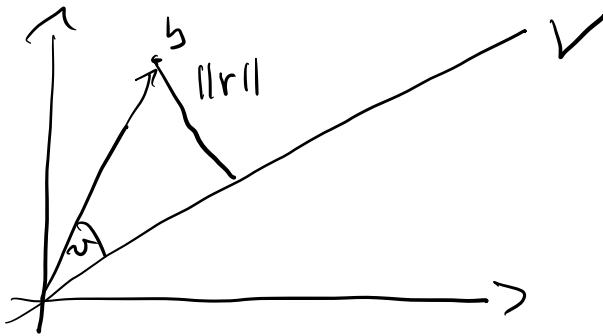
$$= \frac{||b||^2 - ||b - Pb||^2}{||b||^2}$$

$$= 1 - \sin^2(v) = \cos^2(v)$$

$$k_{rel}(b) = \frac{1}{\cos(v)} \, ||P||_2$$

$$r = b - Ax$$

$$\sin(\theta) = \frac{\|r\|_2}{\|b\|_2}$$

# Linear least-squares problems

Conditioning

Now for the least-squares problem $\|\boldsymbol{Ax} - \boldsymbol{b}\|_2$. The relative condition number $\kappa$ in the Euclidean norm is bounded by

▶ With respect to perturbations in $\boldsymbol{b}$:

$$\kappa \leq \frac{\kappa_2(A)}{\cos(\theta)}$$

▶ With respect to perturbations in $\boldsymbol{A}$:

$$\kappa \leq \kappa_2(A) + \kappa_2(A)^2 \tan(\theta)$$

Small residual problems, small angle $\theta$ $\cos(\theta) \approx 1$, $\tan(\theta) \approx 0$: behavior similar to linear system.
Large residual problems, large angle $\theta$ $\cos(\theta) \ll 1$, $\tan(\theta) \approx 1$: behavior very different from linear system because $\kappa_2(A)^2$ shows up

# How should we solve least-squares problems numerically?

We know from the previous slide that if the residual is large, then the condition $\kappa$ is much larger than $\kappa_2(A)$ (closer to $\kappa_2(A)^2$)

▶ This is a poorly condition problem; however, do we care?

# How should we solve least-squares problems numerically?

We know from the previous slide that if the residual is large, then the condition $\kappa$ is much larger than $\kappa_2(A)$ (closer to $\kappa_2(A)^2$)

- ▶ This is a poorly condition problem; however, do we care?
- ▶ If the residual is large, then our $\boldsymbol{Ax}$ won't explain well the right-hand side $\boldsymbol{b}$
- ▶ This means that "our curve doesn't fit well the data" and we probably should try to find another space in which to search for a solution (another $\boldsymbol{A}$ with a range that better approximates the projected right-hand side $\boldsymbol{b}$)

More relevant is the situation with a small residual and then $\kappa \approx \kappa_2(A)$

▶ Here we have a well condition problem; so what could go wrong?

More relevant is the situation with a small residual and then $\kappa \approx \kappa_2(A)$

▶ Here we have a well condition problem; so what could go wrong?

▶ If we choose a numerical method that solves the normal equations

$$\boldsymbol{A}^T \boldsymbol{A} \boldsymbol{x} = \boldsymbol{A}^T \boldsymbol{b}$$

then our problem becomes the problem of solving the linear system with matrix $\boldsymbol{A}^T \boldsymbol{A}$, which has condition number

$$\kappa_2(\boldsymbol{A}^T \boldsymbol{A}) = \kappa_2(\boldsymbol{A})^2$$

More relevant is the situation with a small residual and then $\kappa \approx \kappa_2(A)$

- ▶ Here we have a well condition problem; so what could go wrong?
- ▶ If we choose a numerical method that solves the normal equations

$$\boldsymbol{A}^T \boldsymbol{A} \boldsymbol{x} = \boldsymbol{A}^T \boldsymbol{b}$$

  then our problem becomes the problem of solving the linear system with matrix $\boldsymbol{A}^T \boldsymbol{A}$, which has condition number

$$\kappa_2(\boldsymbol{A}^T \boldsymbol{A}) = \kappa_2(\boldsymbol{A})^2$$

  ⤳ we are back in the situation of a poorly condition problem ("solving a linear system with $\boldsymbol{A}^T \boldsymbol{A}$") even though our original problem (least-squares problem) is well condition

- ▶ Can we do better and solve the least-squares problem (the problem we are originally interested in) without having to solve a problem with condition that grows with $\kappa_2(\boldsymbol{A})^2$ on the way?

Recall that projecting $\boldsymbol{b}$ onto the column span (range) of $\boldsymbol{A}$ was the key step $\rightsquigarrow$ let's try to find a numerical method that computes an orthonormal basis $\boldsymbol{q}_1, \ldots, \boldsymbol{q}_n$ of the rank-$n$ column span of $\boldsymbol{A}$

$$
\boldsymbol{A} = \begin{bmatrix} \boldsymbol{a}_1 & \cdots & \boldsymbol{a}_n \end{bmatrix} \in \mathbb{R}^{m \times n}, \qquad m \geq n
$$

$$
\underbrace{\begin{bmatrix} \boldsymbol{a}_1 & \cdots & \boldsymbol{a}_n \end{bmatrix}}_{\boldsymbol{A}} = \underbrace{\begin{bmatrix} \boldsymbol{q}_1 & \cdots & \boldsymbol{q}_n \end{bmatrix}}_{\boldsymbol{Q}} \underbrace{\begin{bmatrix} r_{11} & r_{12} & \cdots & r_{1n} \\ & r_{22} & & \vdots \\ & & \ddots & \\ & & & r_{nn} \end{bmatrix}}_{\boldsymbol{R}}
$$

with an invertible matrix $\boldsymbol{R}$ so that

$$
\mathrm{span}(\boldsymbol{a}_1, \ldots, \boldsymbol{a}_k) = \mathrm{span}(\boldsymbol{q}_1, \ldots, \boldsymbol{q}_k), \qquad k = 1, \ldots, n
$$

$$\underbrace{\begin{bmatrix} | & & | \\ \boldsymbol{a}_1 & \cdots & \boldsymbol{a}_n \\ | & & | \end{bmatrix}}_{\boldsymbol{A}} = \underbrace{\begin{bmatrix} | & & | \\ \boldsymbol{q}_1 & \cdots & \boldsymbol{q}_n \\ | & & | \end{bmatrix}}_{\boldsymbol{Q}} \underbrace{\begin{bmatrix} r_{11} & r_{12} & \cdots & r_{1n} \\ & r_{22} & & \cdots \\ & & \ddots & \\ & & & r_{nn} \end{bmatrix}}_{\boldsymbol{R}}$$

$\Downarrow$ leads to system of equations $\Downarrow$

$$\boldsymbol{a}_1 = r_{11}\boldsymbol{q}_1$$
$$\boldsymbol{a}_2 = r_{12}\boldsymbol{q}_1 + r_{22}\boldsymbol{q}_2$$
$$\boldsymbol{a}_3 = r_{13}\boldsymbol{q}_1 + r_{23}\boldsymbol{q}_2 + r_{33}\boldsymbol{q}_3$$
$$\vdots$$
$$\boldsymbol{a}_n = r_{1n}\boldsymbol{q}_1 + r_{2n}\boldsymbol{q}_2 + \cdots + r_{nn}\boldsymbol{q}_n$$

This motivates a process for computing the basis $q_1, \ldots, q_n$

▶ At step $j$, we have $q_1, \ldots, q_{j-1}$ that span $\text{span}(a_1, \ldots, a_{j-1})$
▶ We want to find $q_j$ orthonormal to $q_1, \ldots, q_{j-1}$ so that $q_1, \ldots, q_j$ spans $\text{span}(a_1, \ldots, a_j)$
▶ Thus, set

$$v_j = a_j \ominus (q_1^T a_j)q_1 \ominus (q_2^T a_j)q_2 \ominus \cdots \ominus (q_{j-1}^T a_j)q_{j-1}$$

and normalize

$$q_j = \frac{v_j}{\|v_j\|_2}$$

Notice that at step $j$, the quantities $q_1^T a_j, q_2^T a_j, \ldots, q_{j-1}^T a_j$ are the values $r_{j,1}, \ldots, r_{j,j-1}$ and $r_{jj}$ is responsible for the normalization and set to

$$r_{jj} = \|a_j - \sum_{i=1}^{j-1} r_{ij} q_i\|_2$$

This process is the *classical Gram-Schmidt* procedure to compute the $QR$ factorization; however, this process is numerically unstable!