# Solving linear systems of equations

3/29/11

⌐⊥

Our general strategy for solving a system $Ax=b$ will be to __transform__ it to an equivalent, but easier to solve problem (or problems). An example of an easier sub-problem is a "triangular" system $Ux=b$ where

$$U = \begin{bmatrix} u_{11} & u_{12} & \cdots & u_{1n} \\ 0 & u_{22} & \cdots & u_{2n} \\ \vdots & O & \ddots & \vdots \\ 0 & 0 & \cdots & 0 \quad u_{nn} \end{bmatrix}$$ is an upper triangular matrix.

Here is an example, illustrating how such systems are easy to solve

$$\begin{bmatrix} 1 & 2 & 2 \\ 0 & -4 & -6 \\ 0 & 0 & -1 \end{bmatrix} \begin{bmatrix} x_1 \\ x_2 \\ x_3 \end{bmatrix} = \begin{bmatrix} 3 \\ -6 \\ 1 \end{bmatrix}$$

replace $x_3$ with $-1$

Solve: $x_3 = -1$

$$\Downarrow$$

$$\begin{bmatrix} 1 & 2 & 0 \\ 0 & -4 & 0 \\ 0 & 0 & 1 \end{bmatrix} \begin{bmatrix} x_1 \\ x_2 \\ x_3 \end{bmatrix} = \begin{bmatrix} 5 \\ -12 \\ -1 \end{bmatrix}$$

replace $x_2$ with 3

Solve: $x_2 = 3$

$$\begin{bmatrix} 1 & 0 & 0 \\ 0 & 1 & 0 \\ 0 & 0 & 1 \end{bmatrix} \begin{bmatrix} x_1 \\ x_2 \\ x_3 \end{bmatrix} = \begin{bmatrix} -1 \\ 3 \\ -1 \end{bmatrix}.$$

We can write this procedure, formally, in pseudocode

## Back-substitution for Upper Triangular System

for $j = n$ to $1$

    if $u_{jj} = 0$ then <u>stop</u>        [matrix is singular]

    $x_j \leftarrow b_j / u_{jj}$         (1)

    for $i = 1$ to $j-1$

        $b_i \leftarrow b_i - u_{ij} x_j$    (2)

    end

end

By counting how many times the loops are executed, we see that $(n-1)$ divisions (line 1) are required, while line 2 is executed $\frac{n(n-1)}{2}$ times. Overall, the cost of back substitution is proportional to $O(n^2)$.

If
$$L = \begin{bmatrix} l_{11} & 0 & 0 & \cdots & 0 \\ l_{21} & l_{22} & 0 & \cdots & 0 \\ l_{31} & l_{32} & l_{33} & \cdots & 0 \\ \vdots & & & & \vdots \\ l_{n1} & \cdots & & \cdots & l_{nn} \end{bmatrix}$$
is a lower triangular

matrix, a similar procedure can be followed to solve $Lx = b$ :

### Forward substitution for $Lx = b$

for $j = 1$ to $n$

    if $l_{jj} = 0$ then stop         [matrix is singular]

    $x_j \leftarrow b_j / l_{jj}$

    for $i = j+1$ to $n$

        $b_i \leftarrow b_i - l_{ij} x_j$

    end

end

Cost $= O(n^2)$ again.

The forward- and backward substitution processes can be used to solve a non-triangular system by virtue of the following factorization property :

__Thm__   If $A$ is an $n \times n$ matrix, it can be (generally) written as a product:

$$A = LU.$$

where $L$ = lower triangular & $U$ = upper triangular.

Furthermore it is possible to construct $L$ such that all diagonal elements $l_{ii} = 1$

__Algorithm__ : LU factorization by Gaussian Elimination

[ _Note_ : This algorithm executes in-place. I.e. the matrix $A$ is _replaced_ by its LU factorization, in compact form]

for $k=1$ to $n-1$

    if $a_{kk} = 0$ then stop

        for $i = k+1$ to $n$

            $a_{ik} \leftarrow a_{ik} / a_{kk}$

        end

        for $j = k+1$ to $n$

            for $i = k+1$ to $n$

                $a_{ij} = a_{ij} - a_{ik} a_{kj}$

            end

        end

end.

$\approx \dfrac{n^3}{3}$ times

Cost $= O(n^3)$.

More specifically, this algorithm produces a factorization $A = LU$, where:

$$L = \begin{bmatrix} 1 & & & & & \\ l_{21} & 1 & & & & \\ l_{31} & l_{32} & 1 & & & \\ l_{41} & l_{42} & l_{43} & 1 & & \\ \vdots & \vdots & \vdots & & \ddots & \\ l_{m1} & l_{n2} & l_{n3} & \cdots & l_{n,n-1} & 1 \end{bmatrix}$$

$$U = \begin{bmatrix} u_{11} & u_{12} & u_{13} & \cdots & & u_{1n} \\ & u_{22} & u_{23} & \cdots & & u_{2n} \\ & & u_{33} & \cdots & & u_{3n} \\ & & & & \ddots & u_{n-1,n} \\ & & & & & u_{nn} \end{bmatrix}$$

After the in-place factorization algorithm completes, $A$ is replaced by the following "compacted" encoding of $L$ & $U$ together:

$$A = \begin{bmatrix} u_{11} & u_{12} & u_{13} & \cdots & & u_{1n} \\ l_{21} & u_{22} & u_{23} & \cdots & & u_{2n} \\ l_{31} & l_{32} & u_{33} & \cdots & & u_{3n} \\ \vdots & \vdots & \vdots & \ddots & & \vdots \\ & & & & u_{n-1,n-1} & u_{n-1,n} \\ l_{n1} & l_{n2} & l_{n3} & \cdots & l_{n,n-1} & u_{nn} \end{bmatrix}$$

Existence/Uniqueness and pivoting

When computing the LU factorization, the algorithm will **halt** if the diagonal element $a_{kk} = 0$. This can be avoided by swapping rows of A prior to computing the LU factorization. This is done to always select the largest $a_{kk}$ from the equations that follow

$$A \begin{bmatrix} 1 & 2 & 5 & -1 \\ 0 & ⓪ & 3 & 1 \\ 0 & 4 & 1 & -2 \\ 0 & -6 & 0 & 3 \end{bmatrix}$$

largest

$$P = \begin{bmatrix} 1 & 0 & 0 & 0 \\ 0 & 0 & 0 & 1 \\ 0 & 0 & 1 & 0 \\ 0 & 1 & 0 & 0 \end{bmatrix}$$

I with swap columns

$$PA = \begin{bmatrix} 1 & 2 & 5 & -1 \\ 0 & -6 & 0 & 3 \\ 0 & 4 & 1 & -2 \\ 0 & 0 & 3 & 1 \end{bmatrix} \checkmark$$

This is pivoting: The pivot $a_{kk}$ is selected to be nonzero. In this process, we can guarantee uniqueness & existence of LU:

**Thm** If $P$ is a permutation matrix such that all pivots in the Gaussian Elimination of $PA$ are non-zero, and $l_{ii} = 1$, then the LU factorization exists and is unique!

$$PA = LU.$$

## Solving $Ax = b$

- Without pivoting:

$$Ax = b$$

$$\underbrace{LUx}_{=y} = b \implies$$

1. Solve $Ly = b$  (F.S).

2. Solve $Ux = y$  (B.S).

$x$ is the solution.

$$[L, U] = \text{lu}(A) \implies \text{Cost } o(n^3)$$

$$y = b \backslash L \implies \text{Cost } o(n^2)$$

$$x = y \backslash L \implies \text{Cost } o(n^2).$$

If have multiple systems $Ax_i = b_i \implies$ cost of $LU$

only once.

− With pivoting :  For $P^T P = I$

$$Ax = b \iff PAx = Pb \iff P^T P Ax = P^T P b$$

$$L\underset{y}{\underbrace{U x}} = Pb$$

$$Ly = Pb$$

$$Ux = y$$

Matlab

$$[L, U, P] = lu(A);  \longrightarrow  \text{Cost } o(n^3)$$

$$z = Pb;  \longrightarrow  \text{Cost } o(n)$$

$$y = L \backslash z;  \longrightarrow  \text{Cost } o(n^2)$$

$$x = U \backslash y;  \longrightarrow  \text{Cost } o(n^2)$$

Final variant : **Full pivoting**

In this case, when we are in the $k$th step of the Gauss-Elimination / LU procedure, we pick the pivot element among the **entire** $(n-k+1) \times (n-k+1)$ lower-rightmost submatrix of $A$.

e.g. $k=2$ and $Ax=b \implies$

$$\begin{bmatrix} 1 & 2 & 5 & -1 \\ 0 & 0 & 3 & 1 \\ 0 & 4 & 1 & -8 \\ 0 & -6 & 0 & 3 \end{bmatrix} \begin{bmatrix} x_1 \\ x_2 \\ x_3 \\ x_4 \end{bmatrix} = \begin{bmatrix} 4 \\ 7 \\ 8 \\ 2 \end{bmatrix}$$

pivot position

Submatrix

largest element.

In this case, we can bring $(-8)$ to the pivot position $a_{22}$ by permuting **both** rows 2-3 AND columns 2-4. Naturally, we will respectively swap rows 2-3 of the RHS, and rows 2-4 of the unknown vector

Thus :
$$\begin{bmatrix} 1 & 2 & 5 & -1 \\ 0 & -8 & 1 & 4 \\ 0 & 1 & 3 & 0 \\ 0 & 3 & 0 & -6 \end{bmatrix} \begin{bmatrix} x_1 \\ x_4 \\ x_3 \\ x_2 \end{bmatrix} = \begin{bmatrix} 4 \\ 8 \\ 7 \\ 2 \end{bmatrix}$$
is the equivalent system !

This process is encoded in the LU factorization using $\underline{\underline{2}}$ permutation matrices $P$ & $Q$ such that $\boxed{PAQ = LU}$

The solution is then computed via :

$$Ax = b \implies PA\underbrace{QQ^T}_{=I}x = Pb \implies (LU)(Q^Tx) = Pb$$

$$\implies L\underbrace{(UQ^Tx)}_{Y} = Pb \implies y \text{ found via lower triangular solve}$$

$$U\underbrace{Q^Tx}_{z} = y \implies z \text{ found via upper triangular solve}$$

$$Q^Tx = z \implies QQ^Tx = Qz \implies \boxed{x = Qz}$$

In Matlab :

```
[L,U,P,Q] = lu(A);     → Cost  O(n³)
c = Pb;                → Cost  O(n)
y = L\c;               → Cost  O(n²)
z = U\y;               → Cost  O(n²)
x = Q*u;               → Cost  O(n).
```

We saw that in order to avoid dividing by zero (or, near-zero) values in the LU decomposition, we may need to resort to partial or full pivoting.

- Partial pivoting permutes rows, such that the pivot element in the k-th iteration is the largest number in the $(n-k+1)$ lower entries of the k-th column. It is written, in the context of LU decomposition, as

$$PA = LU \qquad (P = \text{permutation})$$

- Full pivoting selects the pivot element in the k-th iteration as the largest element of the $(n-k+1) \times (n-k+1)$ lower-rightmost submatrix of A. It operates by permuting rows AND columns and leads to a LU decomposition of.

$$PAQ = LU.$$

However, there are certain categories of matrices for which we can safely use Gauss elimination or LU decomposition without the need for pivoting (i.e. the pivot elements will never be problematically small).

<u>Def</u>  A matrix A is called <u>diagonally dominant</u>

<u>by columns</u> if the magnitude of every diagonal element

is larger than the sum of magnitudes of all other

entries in the same column, i.e.

For every $i = 1, 2, \ldots, n$         $|a_{ii}| > \sum\limits_{j \neq i} |a_{ji}|$

If the diagonal element exceeds in magnitude the sum

of magnitudes of all other elements in its <u>row</u>, i.e.

for $i = 1, 2, \ldots, n$         $|a_{ii}| > \sum\limits_{j \neq i} |a_{ij}|$

then the matrix is called diagonally <u>dominant by rows</u>.

<u>Def</u>

A ˢʸᵐᵐᵉᵗʳⁱˣ matrix $A \in \mathbb{R}^{n \times n}$ is called positive definite (in short:

SPD : for "symmetric positive definite"), if for any $\underline{x} \in \mathbb{R}^n$ $\underline{x} \neq 0$

we have     $\underline{x}^T A \underline{x} > 0$

If for any $\underline{x} \in \mathbb{R}$ $\underline{x} \neq 0$ we have $\underline{x}^T A \underline{x} \geq 0$, the

matrix is called positive semi-definite.

If the respective properties are $x^T A x < 0$ (or $x^T A x \leq 0$) the

<u>Def</u> the k-th <u>leading principal minor</u> of a matrix $A \in \mathbb{R}^{n \times n}$ is the determinant of the top-leftmost $k \times k$ submatrix of A. Thus if we denote this minor by $M_h$:

$$M_1 = |a_{11}| \qquad M_2 = \begin{vmatrix} a_{11} & a_{12} \\ a_{21} & a_{22} \end{vmatrix} \qquad \cdots \qquad M_k = \begin{vmatrix} a_{11} & \cdots & a_{1k} \\ \vdots & & \vdots \\ a_{k1} & \cdots & a_{kk} \end{vmatrix}$$

<u>Thm</u>  If <u>all</u> leading principal minors (i.e. for $k = 1, 2, 3, \ldots, n$), of the symmetric matrix A are positive, then A is positive definite

\* If $M_k < 0$ for $k = $ odd and $M_k > 0$ for $k = $ even, then A is negative definite.

<u>Thm</u> Pivoting is <u>not</u> necessary when A is diagonally dominant by columns, or symmetric and positive- (or negative-) definite.

These "special" classes of matrices (which appear quite often in engineering and applied sciences) not only make LU decomposition more robust, but also open some additional possibilities for solving $A\underline{x} = \underline{b}$.

## Iterative methods for linear systems.

The general idea is similar to the philosophy of iterative methods we saw for nonlinear equations, i.e. we proceed as follows:

→ We write a (matrix) equation

$$\underline{x} = T\underline{x} + \underline{c}$$

in such a way that this equation is equivalent to $Ax = b$.

→ We start with an initial guess $\underline{x}^{(0)}$ for the solution of $Ax = b$

→ We iterate

$$\underline{x}^{(k+1)} = T\underline{x}^{(k)} + \underline{c}$$

→ If properly designed the sequence $x^{(0)}, x^{(1)}, ..., x^{(k)}, ...$ converges to $\underline{x}^*$, which satisfies $\underline{x}^* = T\underline{x}^* + c$ and subsequently $A\underline{x}^* = b$.

# General design approach

We decompose   $A = D - L - U$

diagonal ↑   lower ↑   ← upper triangular

triangular

e.g.

$$A = \begin{bmatrix} a_{11} & a_{12} & \cdots\cdots & a_{1n} \\ a_{21} & a_{22} & \cdots & a_{2n} \\ \vdots & \vdots & \ddots & \vdots \\ a_{n1} & a_{n2} & \cdots & a_{nn} \end{bmatrix}$$

$$D = \begin{bmatrix} a_{11} & & & \\ & a_{22} & & \\ & & \ddots & \\ & & & a_{nn} \end{bmatrix}$$

$$L = \begin{bmatrix} 0 & & & \\ -a_{21} & 0 & & O \\ -a_{31} & -a_{32} & \ddots & \\ -a_{n1} & -a_{n2}\cdots & -a_{n,n-1} & 0 \end{bmatrix}$$

$$U = \begin{bmatrix} 0 & -a_{12} & -a_{13} & & -a_{1n} \\ & \ddots & -a_{23} & & -a_{2n} \\ & & \ddots & & \\ & & & \ddots & -a_{n-1,n} \\ & & & & 0 \end{bmatrix}$$

# The Jacobi method

$$A = D - L - U$$

$$Ax = b$$

$$(D - L - U)x = b$$

$$Dx = (L+U)x + b$$

$$x = \underbrace{D^{-1}(L+U)x}_{T} + \underbrace{D^{-1}b}_{c} \qquad \left(x = Tx + c\right)$$

Iteration:

$$\underline{x}^{(k+1)} = D^{-1}(L+U)\,\underline{x}^{(k)} + D^{-1}\underline{b} \qquad \underline{\underline{or}}$$

$$\boxed{D\underline{x}^{(k+1)} = (L+U)\underline{x}^{(k)} + \underline{b}} \qquad \text{The Jacobi iteration formula.}$$

$\longrightarrow$ Solution: easy, since we need to solve a linear system of equations with diagonal coefficient matrix

$$\begin{bmatrix} d_1 & & & \\ & d_2 & & \\ & & \ddots & \\ & & & d_n \end{bmatrix} \begin{bmatrix} x_1^{(k+1)} \\ \vdots \\ x_n^{(k+1)} \end{bmatrix} = \begin{bmatrix} c_1 \\ \vdots \\ c_n \end{bmatrix} \longrightarrow x_i^{(k+1)} = \frac{1}{d_i}\,c_i$$

→ Convergence : The Jacobi method is guarranteed to converge when A is diagonally dominant by rows

→ Complexity: Each iteration has a cost associated with

  → Solving $Dx^{(k+1)} = c$ ⟹ $n$ divisions

  → Computing $c = (L+U)x^{(k)} + b$ ⟹ as many additions and multiplications as nonzero entries in A

  worst case $= O(n^2)$, but could be $O(n)$ for sparse matrices

→ Stopping criterion: $\| b - Ax^{(k)} \| < \varepsilon$

  or $\| x^{(k+1)} - x^{(k)} \| < \varepsilon$.

There are 3 forms we will see this algorithm, for different purposes:

(i) Matrix form (for proofs) $Dx^{(k+1)} = (L+U)x^{(k)} + b$

(ii)   Algorithm form (without in-place update)

Each row of $Dx^{(k)} = (L+U)x^{(k)} + b$:

$$a_{ii}\, x_i^{(k+1)} = b_i - \sum_{j \neq i} a_{ij}\, x_j^{(k)}$$

$$\Rightarrow \quad x_i^{(k+1)} = \frac{1}{a_{ii}} \left( b_i - \sum_{j \neq i} a_{ij}\, x_j^{(k)} \right)$$

for $k = 1, 2, \ldots, <max>$

    for $i = 1, 2, \ldots, n$

       $x_i^{(k+1)} \leftarrow \dfrac{1}{a_{ii}} \left( b_i - \sum_{j \neq i} a_{ij}\, x_j^{(k)} \right)$

    end

    <check for convergence

end

iii) In-place algorithm

(replaces $\underline{x}$ with a better estimate)

$\underline{x} \leftarrow$ initial guess

for $k = 1, 2, \ldots, <max>$

   for $i = 1, 2, \ldots, n$

$$x_i^{new} \leftarrow \frac{1}{a_{ii}} \left( b_i - \sum_{j \neq i} a_{ij} x_j \right)$$

   end

   $\underline{x} \leftarrow \underline{x}^{new}$

   check for convergence

end

# The Gauss-Seidel method

We again employ the decomposition $A = D - L - U$

$$Ax = b$$

$$(D - L - U)x = b$$

$$(D - L)x = Ux + b$$

At this point, we place $x^{(k+1)}$ on the LHS and $x^{(k)}$ on the RHS

$$\boxed{(D-L)x^{(k+1)} = Ux^{(k)} + b} \qquad (1)$$

The benefit of the Gauss-Seidel method (1) over Jacobi is the improved convergence, which is guaranteed not only for diagonally dominant matrices, but also for symmetric and _positive definite_ matrices

In terms of complexity, each iteration of (1) amounts to solving a lower-triangular system via forward substitution, i.e. incurs a cost $O(h)$, $h = \#$ of nonzero entries in $A$.

Once again, form (1) is useful for proofs, while the pseudo code version is given as:

(Without in-place update)

$$\underline{x}^{(0)} \longleftarrow \text{initial guess}$$

for $k = 0, 1, \ldots, <\text{max}>$

   for $i = 1, 2, \ldots, n$

$$x_i^{(k+1)} \longleftarrow \frac{1}{a_{ii}}\left(b_i - \sum_{j<i} a_{ij} x_j^{(k+1)} - \sum_{j>i} a_{ij} x_j^{(k)}\right)$$

   end

   <check for convergence

end

---

(In-place)

$$\underline{x} \longleftarrow \text{initial guess}$$

for $k = 0, 1, \ldots, <\text{max}>$

   for $i = 1, 2, \ldots, n$

$$x_i \longleftarrow \frac{1}{a_{ii}}\left(b_i - \sum_{i \neq j} a_{ij} x_j\right) \quad (*)$$

   end

   <check for convergence>

end.

Compare: Jacobi

$$\underline{x} \longleftarrow \text{initial guess}$$

for $k = 0, 1, \ldots, <\text{max}>$

   for $i = 1, 2, \ldots, n$

$$x_i^{new} \longleftarrow \frac{1}{a_{ii}}\left(b_i - \sum_{i \neq j} a_{ij} x_j\right) \quad (**)$$

   end

$$\underline{x} \longleftarrow \underline{x}^{new} \quad (***)$$

end

---

The <u>real</u> difference in performance is :

Gauss-Seidel is generally <u>serial</u> in nature (although parallel variants exist), while steps (**) & (***) are highly parallel.
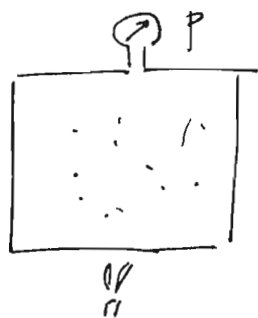
# Overdetermined systems

So far, we considered linear systems $Ax = b$ with the same number of equations & unknowns (i.e. $A \in \mathbb{R}^{n \times n}$). In the case where $A \in \mathbb{R}^{m \times n}$, with $m > n$ (more equations) the existence of a true solution is not guaranteed; in this case we look for the "best possible" substitute for a solution. Before analyzing what that means, let's look at how such problems arise:

e.g. In an experiment, we measure the pressure of a gas in a closed container, as a function of the temperature



From physics:
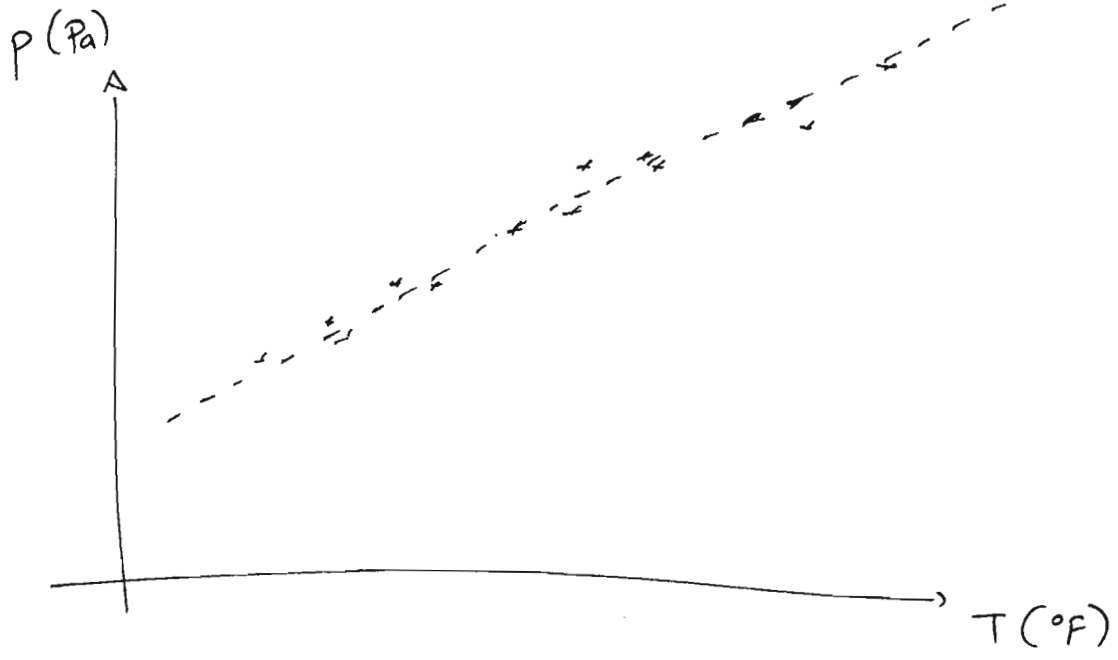
$$pV = nR \frac{5}{9} (T + 459.67)$$

$$T: \text{Fahrenheit}.$$

i.e. $p = \alpha T + \beta$

$$\left( \alpha = \frac{5nR}{9V}, \quad \beta = \frac{5nR \cdot 459.67}{9V} \right)$$
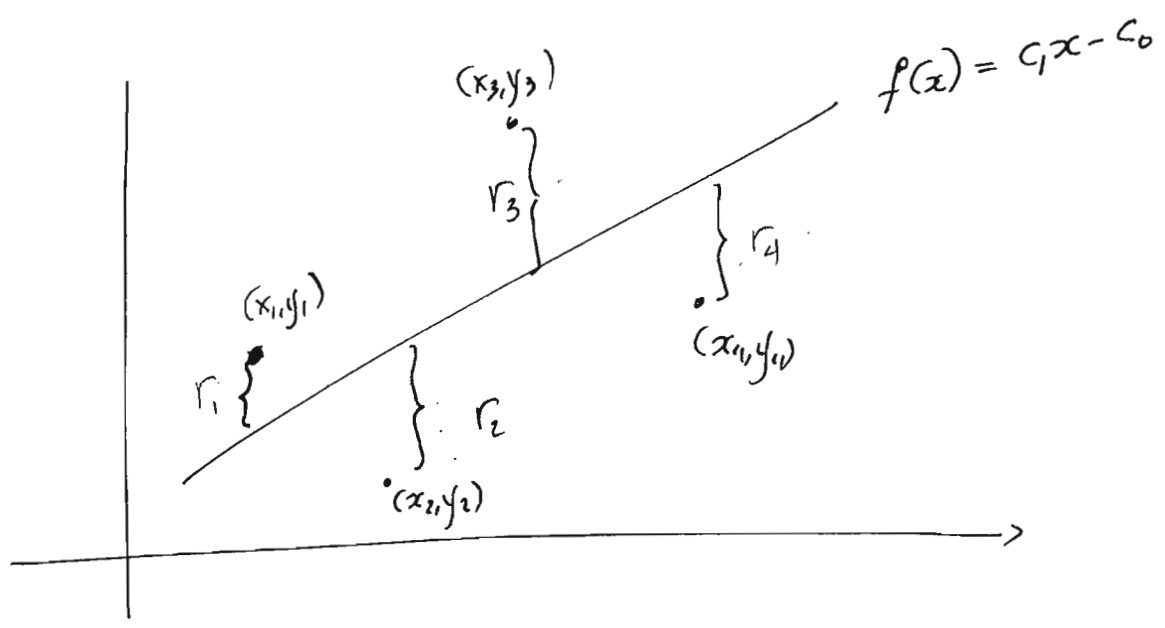
What are $\alpha$ & $\beta$ ?

experimentally :

$p$ (Pa)



$T$ (°F)

The measurements should ideally lie on a straight

line $y = c_1 x + c_0$, but do not, due to measurement

error. If we have $\underline{n}$ measurement pairs $(x_1, y_1), \ldots, (x_n, y_n)$

we would have wanted:

$$
\left.
\begin{aligned}
y_1 &= c_1 x_1 + c_0 \\
y_2 &= c_1 x_2 + c_0 \\
&\vdots \\
y_n &= c_1 x_n + c_0
\end{aligned}
\right\}
\Rightarrow
\underbrace{\begin{bmatrix} x_1 & 1 \\ x_2 & 1 \\ \vdots & \\ x_n & 1 \end{bmatrix}}_{A}
\underbrace{\begin{bmatrix} c_1 \\ c_0 \end{bmatrix}}_{x}
=
\underbrace{\begin{bmatrix} y_1 \\ \vdots \\ \\ y_n \end{bmatrix}}_{b}
$$

Here $\underset{\substack{\uparrow \\ n \times 2}}{A} \underset{\substack{\uparrow \\ 2 \times 1}}{x} = \underset{\substack{\uparrow \\ n \times 1}}{b}$ is a rectangular system.

We cannot hope to find a true solution to this system. Instead let's try to find an "approximate" solution, such that $Ax \approx b$.

Let's look at the residual of this "interpolation"



the residual of the approximation at each data point is

$$r_i = y_i - f(x_i) = y_i - c_1 x_i - c_0$$

If we write the vector of all residuals:

$$\underline{r} = \begin{bmatrix} r_1 \\ \vdots \\ r_n \end{bmatrix} = \begin{bmatrix} y_1 - c_1 x_1 - c_0 \\ y_2 - c_1 x_2 - c_0 \\ \vdots \\ y_n - c_1 x_n - c_0 \end{bmatrix} = \begin{bmatrix} y_1 \\ \vdots \\ y_n \end{bmatrix} - \begin{bmatrix} x_1 & 1 \\ x_2 & 1 \\ \vdots \\ x_n & 1 \end{bmatrix} \begin{bmatrix} c_1 \\ c_0 \end{bmatrix} = \underline{b} - A\underline{x}$$

Although we can't find an $\underline{x}$ s.t. $A\underline{x} = \underline{b}$ (thus $\underline{r} = 0$) we can at least try to make $\underline{r}$ __small__.
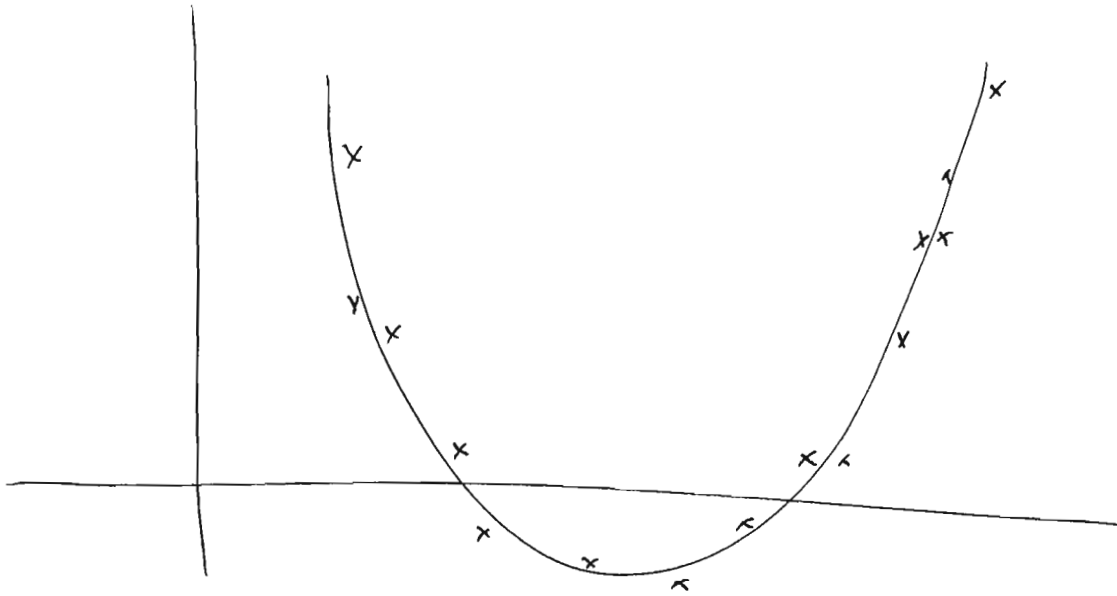
# Another example

Find the best parabola $f(x) = c_2 x^2 + c_1 x + c_0$ that fits measurements $(x_1, y_1), \ldots (x_n, y_n)$



We would like

$$f(x_1) \approx y_1$$
$$f(x_2) \approx y_2$$
$$\vdots$$
$$f(x_n) \approx y_n$$

$$c_2 x_1^2 + c_1 x_1 + c_0 \approx y_1$$
$$c_2 x_2^2 + c_1 x_2 + c_0 \approx y_2$$
$$\vdots$$
$$c_2 x_n^2 + c_1 x_n + c_0 \approx y_n$$

$$\underbrace{\begin{bmatrix} x_1^2 & x_1 & 1 \\ x_2^2 & x_2 & 1 \\ & \vdots & \\ x_n^2 & x_n & 1 \end{bmatrix}}_{A} \underbrace{\begin{bmatrix} c_2 \\ c_1 \\ c_0 \end{bmatrix}}_{x} \approx \underbrace{\begin{bmatrix} y_1 \\ \vdots \\ y_n \end{bmatrix}}_{b}$$

Once again we would like to make $\underline{r} = \underline{b} - A\underline{x}$ as small as possible.

How do we quantify $r$ being small?

$\Rightarrow$ Using a norm!

We could ask that $\longrightarrow$ $\|r\|_1$

$\|r\|_2$ be as small as possibl[e]

or $\|r\|_\infty$

Any of these norms would be intuitive to consider for minimization (esp. $1$- and $\infty$-norms are very intuitive).

However, we typically use the $2$-norm for this purpose, because it's the easiest to work with in this problem!

<u>Def</u> The <u>least squares solution</u> of the overdetermined system $A\underline{x} \approx \underline{b}$ is the vector $\underline{x}$ that minimizes

$$\|\underline{r}\|_2 = \|\underline{b} - A\underline{x}\|_2$$

Define    $Q(\underline{x}) = Q(x_1, x_2, \ldots, x_n)$

$= \| \underline{b} - A\underline{x} \|_2^2$    where    $\underline{x} = (x_1, \ldots, x_n)$

and $A \in \mathbb{R}^{m \times n}$   $b \in \mathbb{R}^m$    $(m > n)$. The least squares

solution is the set of values $x_1, \ldots, x_n$ that $\underline{minimize}$

$Q(x_1, \ldots, x_n)$ !

$Q(x_1, \ldots, x_n) = \| \underline{b} - A\underline{x} \|_2^2 = \| \underline{r} \|_2^2 = \sum_{j=1}^{m} r_j^2$

$\underline{r} = \underline{b} - A\underline{x}$   $\Rightarrow r_i = b_i - [Ax]_i$  $\Rightarrow r_i = b_i - \sum a_{ij} x_j$

$\Rightarrow Q(x_1, \ldots, x_n) = \sum_{i=1}^{m} \left( b_i - \sum_{j=1}^{n} a_{ij} x_j \right)^2$.

If $x_1, \ldots, x_n$ are those that $\underline{minimize}$ $Q$, then:

$\dfrac{\partial Q}{\partial x_1} = 0$

$\dfrac{\partial Q}{\partial x_2} = \ldots = 0$     } in order to guarantee

$\vdots$                a minimum

$\dfrac{\partial Q}{\partial x_n} = 0$.

$$\frac{\partial Q}{\partial x_k} = \frac{\partial}{\partial x_k} \left( \sum_{i=1}^{m} \left( b_i - \sum_{j=1}^{n} a_{ij} x_j \right)^2 \right)$$

$$= \sum_{i=1}^{m} \frac{\partial}{\partial x_k} \left( b_i - \sum_{j=1}^{n} a_{ij} x_j \right)^2$$

$$= \sum_{i=1}^{m} 2 \underbrace{\left( b_i - \sum_{j=1}^{n} a_{ij} x_j \right)}_{r_i} \frac{\partial}{\partial x_n} \left( b_i - \sum_{j=1}^{n} a_{ij} x_j \right)$$

$$= \sum_{i=1}^{m} r_i a_{ik} = \sum_{i=1}^{m} [A^T]_{ki} \, r_i = [A^T r]_k$$

Thus:

$$\frac{\partial Q}{\partial x_1} = 0 \implies [A^T r]_1 = 0$$

$$\vdots$$

$$\frac{\partial Q}{\partial x_n} = 0 \implies [A^T r]_n = 0$$

$$\boxed{A^T \underline{r} = \underline{o}}$$

Since $\underline{r} = \underline{b} - A\underline{x}$ , we have:

$$0 = A^T r = A^T(b - Ax) = A^T b - A^T A x \implies$$

$$\implies \boxed{A^T A \underline{x} = A^T b}$$

The system above is called the <u>normal equations system</u>
it is a <u>square</u> system, that has as solution the
least-squares approximation of $Ax \approx b$

$$\underbrace{A^T \underbrace{A}_{} \underbrace{x}_{}}_{} = \underbrace{A^T \underbrace{b}_{}}_{}$$

$$\underset{n\times m}{} \underset{m\times n}{} \underset{n\times 1}{} \quad \underset{n\times m}{} \underset{m\times 1}{}$$

$$\underset{n\times n}{} \quad \underset{n\times 1}{} \qquad \underset{n\times 1}{}$$

The normal equations <u>always</u> have a solution (with the
simple condition that the columns of $A$ have to be
linearly independent (usually true) ).

<u>Problem</u> : The condition number of $A^T A$ is the
<u>square</u> of that of $A$ (if $A$ was square itself !).

## Overdetermined systems

We examined the case of systems $Ax = b$ where $A \in \mathbb{R}^{m \times n}$ and $\underline{m > n}$. In general, a true solution does not exist. We define however the <u>least squares solution</u> as the vector $\underline{x}$ that minimizes

$$\underline{x} = \arg\min \| \underline{b} - A\underline{x} \|_2^2$$

$\underline{x}$ is given by the solution to the <u>system of normal equations</u>

$$A^T A \underline{x} = A^T \underline{b} \qquad (L)$$

System (1) is square ($n \times n$) and invertible (if $A$ has linearly independent columns). However, the condition number of $A^T A$ could be very poor... for example, if $A$ was square, we would have $\text{cond}(A^T A) = [\text{cond}(A)]^2$.

An alternative method that does not suffer from this problematic conditioning is the $QR$ factorization.

<u>Def</u> An $n \times n$ matrix $Q$ is called <u>orthogonal</u> iff

$$Q^T Q = Q Q^T = I.$$

<u>Thm</u> Let $A \in \mathbb{R}^{m \times n}$ $(m > n)$ have linearly independent columns. Then a decomposition

$$A = QR$$

exists, such that $Q \in \mathbb{R}^{m \times m}$ is orthogonal and

$R \in \mathbb{R}^{m \times n}$ is upper triangular (i.e. $R = \begin{bmatrix} r_{11} & r_{12} & \cdots & r_{1n} \\ & r_{22} & \cdots & r_{2n} \\ & O & \ddots & \\ O & O & \cdots & r_{nn} \\ & & & O \\ O & O & \cdots & O \end{bmatrix}$ )

Additionally, given that $A$ has linearly independent columns, we guarantee that $r_{ii} \neq 0$.

In Matlab, the QR decomposition is obtained via the $qr$ function, i.e

$$[Q, R] = qr(A);$$

Now, let us write

$$Q = \left[ \hat{Q} \mid Q^* \right]$$

where $\hat{Q} \in \mathbb{R}^{m \times n}$ contains the first $m$ columns of $Q$

and $Q^* \in \mathbb{R}^{(m-n) \times m}$ contains the last $(m-n)$ columns

Respectively, we write:

$$R = \left[\begin{array}{c} \hat{R} \\ \hline O_{(m-n) \times n} \end{array}\right] \quad$$ where $\hat{R} \in \mathbb{R}^{n \times n}$ (and upper triangular) contains the first $n$ rows of $R$.

$\hat{R}$ is also <u>nonsingular</u> (for lin.ind. columns of

We can verify the following:

→ $\hat{Q}^T \hat{Q} = I_n$ (although $\hat{Q}\hat{Q}^T \neq I_m$ !)

Proof: $[\hat{Q}^T \hat{Q}]_{ij} = \displaystyle\sum_{k=1}^{m} [\hat{Q}^T]_{ik} [\hat{Q}]_{kj}$
$1 \leq i,j \leq n$

$$= \sum_{k=1}^{m} [\hat{Q}]_{ki}[\hat{Q}]_{kj} = \sum_{k=1}^{m} [Q]_{ki}[Q]_{kj}$$

$$= [Q^T Q]_{ij} = [I_m]_{ij}$$

→ $\underset{(m \times n)}{A} = \underset{(m \times m)}{Q}\underset{(m \times n)}{R} = \underset{(m \times n)}{\hat{Q}} \cdot \underset{(n \times n)}{\hat{R}}$ .

Proof : Similar.

The factorization $A = \hat{Q} \cdot \hat{R}$ is the so-called economy-size QR factorization, and computed in Matlab as:

$$[\hat{Q}, \hat{R}] = qr(A, 0);$$

Once we have $\hat{Q}$ & $\hat{R}$ computed,
we observe that the normal equations can
be written as:

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

using $A = \hat{Q}\hat{R}$ $\Longrightarrow$

$$\hat{R}^T \hat{Q}^T \underbrace{\hat{Q} \hat{R}}_{= I_n} \underline{x} = \hat{R}^T \hat{Q}^T b$$

$= I_n$

$$\hat{R}^T \hat{R} \underline{x} = \hat{R}^T \hat{Q}^T \underline{b}$$

$\hat{R}$ is invertible
$\Longrightarrow$

$$(\hat{R}^{-T})(\hat{R}^T \hat{R} x) = (\hat{R}^{-T})(\hat{R}^T \hat{Q}^T \underline{b})$$

$$\Longrightarrow \boxed{\hat{R}\underline{x} = \hat{Q}^T \underline{b}} \quad (*)$$

Least squares solution
using QR decomposition

Matlab:    $[\hat{Q}, \hat{R}] = qr(A, 0)$

$$z = \hat{Q}' * b;$$

$$x = \hat{R} \backslash z;$$

<u>Benefit</u>: We can show that $\text{cond}(A^T A) = [\text{cond}(\hat{R})]^2$,

thus, equation (*) is <u>much better</u> conditioned than the
normal equations system.