Recall from last lecture:

Theorem (Projection Theorem):

Let $X \leq \mathbb{R}^n$ be a subspace of $\mathbb{R}^n$.
Then for any $y \in \mathbb{R}^n$, there exists a
unique $x^* \in X$ such that

$$\|y - x^*\| \leq \|y - x\| \quad \forall x \in X.$$

The point $x^* \in X$ is uniquely characterized
by $(y - x^*) \in X^\perp$. We call $x^*$ the projection
of $y$ onto $X$.

explicit formula:

$$\circ \quad x^* = \underbrace{A(A^TA)^{-1}A^T}_{} \, y$$

$$= AA^\dagger =: P_X$$

viewed as a constrained optimization problem:

$\circ \quad x^*$ is optimizer of

$$\min_x \|y - x\|$$
$$\text{s.t.} \quad x \in X \leq \mathbb{R}^n$$

Today:

Theorem: Consider an underdetermined system $Ax = b$
where $A \in \mathbb{R}^{m \times n}$ has linearly independent ROWS
Then the unique minimum norm solution
is characterized by $x^* \in \text{range}(A^T)$.
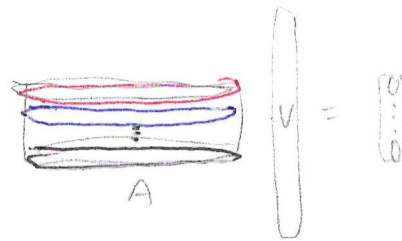Using this information alone, the explicit formula

$$x^* = A^T(AA^T)^{-1} b.$$

can be derived

# First, observe:

$A \in \mathbb{R}^{m \times n}$, $m \leq n$

$$\text{null}(A) = \{ v \in \mathbb{R}^n : Av = 0 \}$$

null(A) is set of vectors $v$ that are
orthogonal to all rows of A
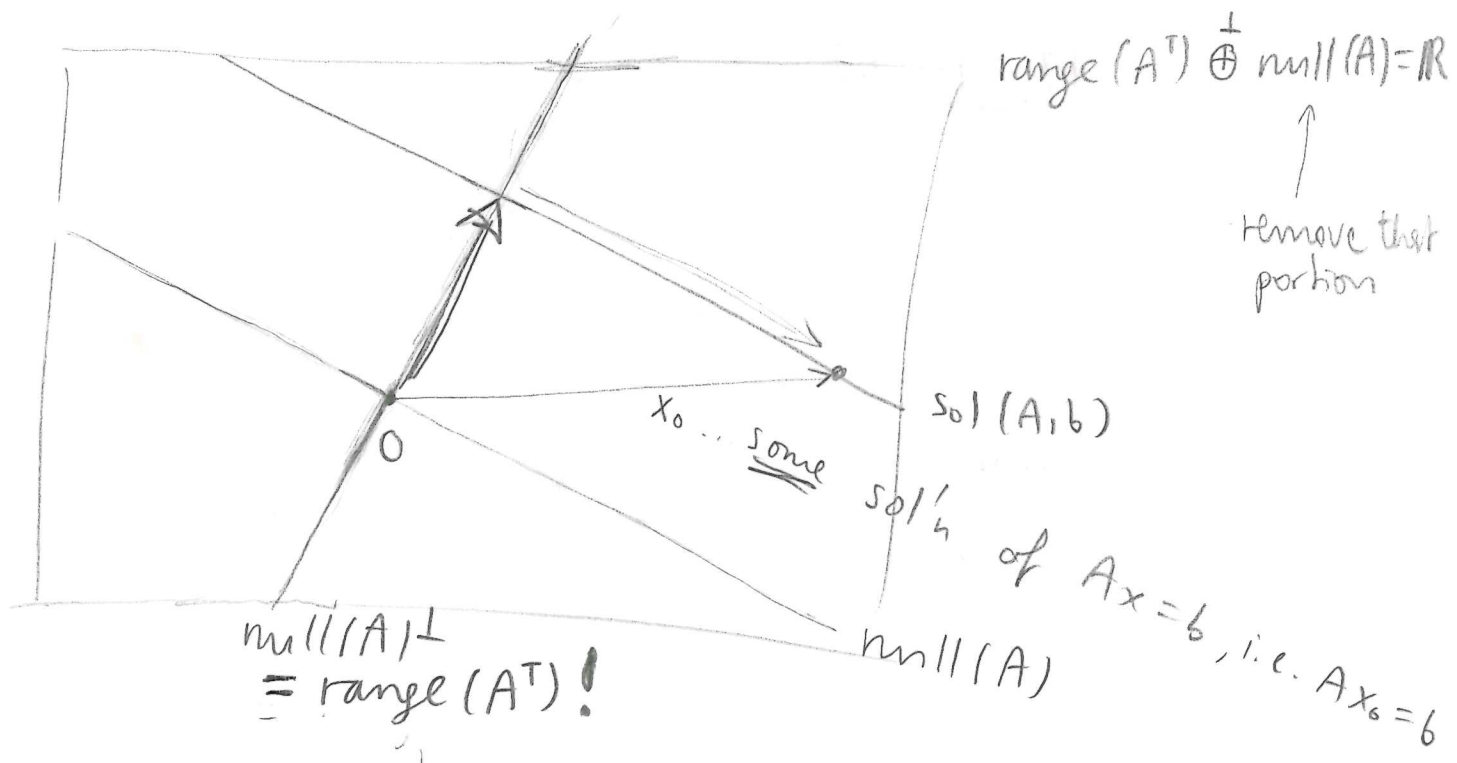$$= \text{columns of } A^T. \quad (\text{by def. of } A^T)$$

$$\begin{bmatrix} a_{11} \\ \vdots \\ u_m \end{bmatrix} = \begin{bmatrix} \\ \end{bmatrix} u_1 + \ldots + \begin{bmatrix} \\ \end{bmatrix} u_m \in \text{range}(A^T)$$

$\Rightarrow$ vectors in null(A) are orthogonal to
vectors in range($A^T$).

$$\text{null}(A)^{\perp} = \text{range}(A^T)$$
$$\left( \text{range}(A^T)^{\perp} = \text{null}(A) \right)$$

range$(A^T) \overset{\perp}{\oplus}$ null$(A) = \mathbb{R}$

↑

remove that portion



0

null$(A)^{\perp}$
$=$ range$(A^T)$ !

$x_0$ ... some sol'n of $Ax = b$, i.e. $Ax_0 = b$

sol$(A, b)$

null$(A)$

## PROOF #1 : (proof by picture)

Project $x_0$ onto range$(A^T)$ by using the formula

(*) $P_X = A(A^TA)^{-1}A^T$ for projecting onto range$(A)$.

That is, replace all appearances of $A$ w/ $A^T$ in (*) and vice versa.

$x^*$ ... optimal sol of $Ax = b$

$= A^T(AA^T)^{-1}\underbrace{Ax_0}_{= b}$

← sol'n now stated in terms of $b$ and independent of choice of particular solution ! ☺

$$\boxed{x^* = A^T(AA^T)^{-1}b}$$

# PROOF #2 (more formal, using the Projection Theorem)

Recall that for an underdetermined system $Ax = b$, once we have A SOLUTION (one of infinitely many) $x_0$, we can generate all the other solutions by adding elements of **null(A)** to that solution $x_0$.

$$A(x_0 + h) = Ax_0 + Ah = b \Rightarrow x_0 + h \text{ is sol'n also.}$$

where $Ax_0 = b$, $Ah = 0$, $h$ in null(A)

Denote $x^*$ the min. norm sol'n of $Ax = b$ (to be sought)

Write $x^* = x_0 - y^*$, with $y^* \in$ null(A).

($x_0$ some "anchor")

Since $x^*$ is min. norm sol'n, it has to hold that

$$\| x_0 - y^* \| \leq \| x_0 - y \| \quad \forall y \in \text{null}(A) =: X$$

This fits description of <u>Projection Theorem</u> precisely.

Thus $y^* \in$ null **(A)** is <u>uniquely</u> characterized by

$$x^* = x_0 - y^* \in \text{null}(A)^{\perp} = \underline{\text{range}(A^T)}$$

Now we know that the ~~optimal~~ $x^* \in$ range$(A^T)$, we **substitute**

$$x^* = A^T u^*$$

and plug it into $Ax = b$ as follows:

($Ax$ is $A^T u$)

$$(AA^T) u = b \Rightarrow u = (AA^T)^{-1} b \Rightarrow \underline{x^* = A^T(AA^T)^{-1}b}$$

(u unique)

A having linearly independent rows means full row rank and thus $AA^T$ (square matrix) is exactly invertible!

(back of envelope proof portion)

# Conjugate Gradient Method

Particularly useful for solving quadratic problems

$$\min_{x \in \mathbb{R}^n} \quad f(x) = \frac{1}{2} x^T Q x - b^T x \quad , \quad Q > 0,$$

or equivalently, the linear system $Qx = b$
since $(\nabla f)(x) = 0 \iff Qx - b = 0$.

For large $n$, $Qx = b$ can be surprisingly hard
to solve, so having an iterative method is desired.

"Standard" Gradient Descent with exact line
search produces directions that are necessarily

(HW3/ P1) <u>orthogonal</u>, which leads to a zig-zag path
that takes unnecessarily long to approach
$x^*$, especially towards later iterations.

The Conjugate Gradient Method guarantees
convergence to $x^*$ in $n$ steps!

**Def.** Given an $n \times n$ symmetric matrix $Q$, we call a set of $n$ non-zero vectors $\{d_1, \ldots, d_n\} \subset \mathbb{R}^n$ $\underline{Q\text{-conjugate}}$ if

$$\underbrace{d_i^T Q \, d_j}_{= \langle d_i, d_j \rangle_Q} = 0 \qquad \forall_{i \neq j} \quad i, j \in \{1, \ldots, n\}.$$

inner product induced by $Q$.

**Proposition:** If $Q > 0$ and $d_1, \ldots, d_n$ are $Q$-conjugate, then $d_1, \ldots, d_n$ are linearly independent!

**Proof:** HW 6, P4.

# Main idea of Conjugate Gradient Method:

With $d_0, d_1, \ldots, d_{n-1}$ Q-conjugate directions,
(by the proposition, this is a basis of $\mathbb{R}^n$!),

The solution $x^*$ of the quadratic optimization problem can be expressed as

$$x^* = \sum_{i=0}^{n-1} \alpha_i d_i$$

$\uparrow$ for some suitable $\alpha_i$

Applying $d_i^T Q$ from the left onto $x^* = \sum_{i=0}^{n-1} \alpha_i d_i$

yields
$$d_i^T Q x^* = \alpha_i d_i^T Q d_i \qquad (d_i^T Q d_j = 0 \text{ if } i \neq j)$$

$$\implies \alpha_i = \frac{d_i^T Q x^*}{d_i^T Q d_i} = \frac{d_i^T b}{d_i^T Q d_i} \leftarrow \text{now independent of } x^*! \text{ NICE}$$

$$Q x^* = b$$

Thus
$$x^* = \sum_{i=0}^{n-1} \alpha_i d_i = \sum_{i=0}^{n-1} \frac{d_i^T b}{d_i^T Q d_i} d_i.$$

The above expansion for $x^*$ can be considered to be the result of an iterative process of $n$ steps where the $i$th step adds "$+ \alpha_i d_i$":

$$\underline{i=0} : \quad \alpha_0 d_0$$
$$\underline{i=1} : \quad \alpha_0 d_0 + \alpha_1 d_1$$
$$\vdots$$
$$\underline{i=n-1} : \quad x^* = \alpha_0 d_0 + \alpha_1 d_1 + \ldots + \alpha_{n-1} d_{n-1}$$

The following result generalizes the above observations by incorporating an initialization $x^0 \neq 0$. The derivation is analogous and left as an exercise (HW 6 / P5).

## Theorem (Conjugate Direction Theorem)

Let $d_0, d_1, \ldots, d_{n-1}$ be a set of nonzero $Q$-conjugate vectors ($Q > 0$). For any $x^0 \in \mathbb{R}^n$, the sequence generated via

$$x^{k+1} = x^k + \alpha_k d_k$$

$$\uparrow$$

$$\alpha_k = - \frac{g_k^T d_k}{d_k^T Q d_k} \quad \text{where } g_k = Q x^k - b,$$

$$\left( \begin{array}{c} \text{gradient of} \\ f(x) = \frac{1}{2} x^T Q x - b^T x \end{array} \right)$$

converges to the unique solution $x^*$ of $Qx = b$ after $n$ steps, i.e. $x^n = x^*$.