In general, it is <u>not</u> a good idea to solve the normal eqn:

$$A^T A \, \mathbb{x} = A^T \, \mathbb{b}$$

by explicitly forming $A^T A$, and then compute $(A^T A)^{-1}$.

<u>why?</u>

1) Forming $A^T A \to$ loss of info.

2) $\kappa(A^T A) = \kappa(A)^2$, i.e.,

the cond. number of $A^T A$ is much worse than that of $A$ in general.

This example is a bit extreme...   Show previous MATLAB example

<u>EX</u>.  <u>Forming $A^T A$ is bad</u>.

$$A = \begin{bmatrix} 1 & 1 \\ \varepsilon & 0 \\ 0 & \varepsilon \end{bmatrix}, \quad \text{say } \varepsilon = 10^{-8}$$

in double precision floating point sys.

Then $A^T A = \begin{bmatrix} 1+\varepsilon^2 & 1 \\ 1 & 1+\varepsilon^2 \end{bmatrix}$

$$\approx \begin{bmatrix} 1 & 1 \\ 1 & 1 \end{bmatrix} \quad \begin{array}{l} \text{because} \\ \varepsilon^2 = 10^{-16} \end{array}$$

<u>How about the condition numbers?</u>

$\kappa(A) \approx 1.4142 \times 10^8$  already bad.

$\kappa(A^T A) \approx +\infty$ in double precision.

If we set $\varepsilon = 10^{-7}$ instead of $10^{-8}$,
then $\kappa(A) \approx 1.4142 \times 10^7$
$$\kappa(A^T A) \approx 1.9903 \times 10^{14}$$
This is still too bad to get any
reliable LS solution for such $A$.

Often such situations occur
when some of the column vectors
of $A$ are "close to parallel", i.e.,
they become almost linearly dependent.

<u>Def.</u> Let $A \in \mathbb{R}^{m \times n}$. Then
$A$ is called rank deficient if
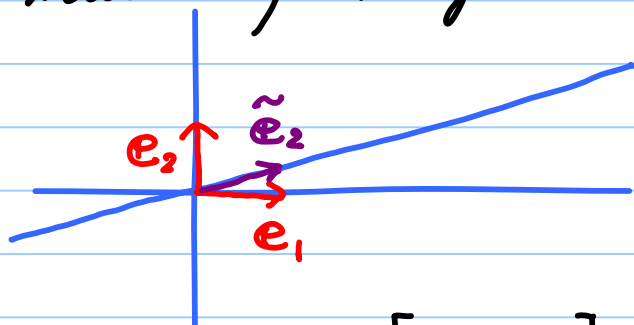$\mathrm{rank}(A) < \min(m, n)$.

i.e., if $A$ is not of full rank.

In general, we should avoid
computing a solution for a given
LS problem by forming $A^T A$ explicitly
and computing $(A^T A)^{-1} A^T b$.
$\Rightarrow$ Better to use the methods
based on QR decomposition or
SVD (we'll discuss these later
in this course.)

# Orthogonality

The above discussion should convince you that $A$ is quite "_good_" if its column vectors are mutually orthogonal.



Suppose $A = [e_1 \ e_2]$, $\tilde{A} = [e_1 \ \tilde{e}_2]$ in $\mathbb{R}^2$. You can see that $A$ is much more "well-balanced" and convenient than $\tilde{A}$. For example, suppose we want to represent $X = [1, 1]^T$ in the basis of $\{e_1, e_2\}$ and that of $\{e_1, \tilde{e}_2\}$. Then the coefficient of $X$ w.r.t. $\{e_1, e_2\}$ is the same as $X$ itself since $A^{-1}X = AX = X$

$$A = I \text{ in } \mathbb{R}^2$$

But $\tilde{A}^{-1}X$ behaves badly.

**Why?** Say $C = \tilde{A}^{-1}X$, $C = [c_1, c_2]^T$

Then $X = \tilde{A}C = [e_1 \ \tilde{e}_2]\begin{bmatrix} c_1 \\ c_2 \end{bmatrix}$

$$= c_1 e_1 + c_2 \tilde{e}_2$$

But $X = e_1 + e_2$, i.e.,

$$e_1 + e_2 = c_1 e_1 + c_2 \tilde{e}_2$$

Taking an inner product with $e_2$ on both sides yields

$$\underbrace{e_2^T(e_1 + e_2)}_{\substack{\| \\ e_2^T e_2 \\ \| \\ \|e_2\|_2^2 = 1.}} = \underbrace{e_2^T(c_1 e_1 + c_2 \tilde{e}_2)}_{\substack{\| \\ c_1 \underbrace{e_2^T e_1}_{=0} + c_2 e_2^T \tilde{e}_2 \\ \| \\ c_2 e_2^T \tilde{e}_2}}$$

$$\Rightarrow \quad 1 = c_2 e_2^T \tilde{e}_2$$

$$\Rightarrow \quad c_2 = \frac{1}{e_2^T \tilde{e}_2}$$

<span style="color:red">could be huge if $\tilde{e}_2$ is close to perpendicular to $e_2$, i.e., close to parallel to $e_1$ !!</span>

✱ <u>Orthogonal Vectors</u>

<u>Def.</u> • Two vectors $x, y \in \mathbb{R}^m$ are said to be <u>orthogonal</u> if $x^T y = 0$. <span style="color:green">So, the zero vector $0$ is orthogonal to <u>any</u> vector.</span>

• Two <u>sets</u> of vectors $X, Y$ are said to be <u>orthogonal</u> if $\forall x \in X, \forall y \in Y, \quad x^T y = 0$.

• A set of vectors $S$ is said to be <u>orthogonal</u> if $\forall x \in S, \forall y \in S, x \neq y$ $x^T y = 0$.

- A set of vectors $S$ is said to be <u>orthonormal</u> if $S$ is orthogonal and $\forall x \in S$, $\|x\|_2 = 1$.

<span style="color:red">even more balanced!</span>

<u>Thm</u> The vectors in an <span style="color:red">orthogonal</span> set $S$ are <span style="color:red">linearly independent</span>.

(Proof) Let $S = \{v_1, \cdots, v_n\}$
Suppose they are <u>not</u> lin. indep.
Then $\exists \, v_k \in S$ s.t. $v_k \neq 0$ and

$$v_k = \sum_{\substack{i=1 \\ i \neq k}}^{n} c_i v_i \quad \text{with} \quad c \neq 0$$

$$c = [c_1, \cdots, c_{k-1}, c_{k+1}, \cdots, c_n]^T$$

Since $S$ is an orthogonal set,
$$v_j^T v_i = 0 \quad \text{for} \quad \forall j \neq i.$$
But $\quad v_k^T \left( \sum_{\substack{i=1 \\ i \neq k}}^{n} c_i v_i \right) = \sum_{\substack{i=1 \\ i \neq k}}^{n} c_i \underbrace{v_k^T v_i}_{\color{red}0} = 0$

$\Leftrightarrow v_k^T v_k = 0$
$\Leftrightarrow \|v_k\|^2 = 0 \Leftrightarrow v_k = 0$ # contradiction!

★ <u>Components of a vector</u>

<span style="color:green">SLOGAN</span> " Inner products can be used to decompose arbitrary vectors into orthogonal components ! "

Suppose $\{ q_1, \cdots, q_n \} \subset \mathbb{R}^m$ is an orthonormal set. $\quad q_j \in \mathbb{R}^m, \ 1 \le j \le n.$

Let $v$ be an arbitrary vector in $\mathbb{R}^m$.

$$r = v - (q_1^T v) q_1 - (q_2^T v) q_2 - \cdots - (q_n^T v) q_n$$

$\uparrow$ residual vector is $\perp$ to $\{ q_1, \cdots, q_n \}$

<u>why?</u>

$$q_j^T r = q_j^T v - (q_1^T v)\underbrace{q_j^T q_1}_{=0} - \cdots - (q_{j-1}^T v) \underbrace{q_j^T q_{j-1}}_{=0}$$
$$- (q_j^T v)\underbrace{q_j^T q_j}_{1} - (q_{j+1}^T v)\underbrace{q_j^T q_{j+1}}_{=0} - \cdots - (q_n^T v)\underbrace{q_j^T q_n}_{=0}$$

$$= q_j^T v - q_j^T v = 0$$

This is true for any $j = 1, \cdots, n$

$$\Rightarrow \quad v = r + \sum_{i=1}^{n} (q_i^T v) q_i$$

any vector in $\mathbb{R}^m$ $\nearrow$

$$= r + \sum_{i=1}^{n} (q_i q_i^T) v$$

$$\perp$$

$$= r + Q Q^T v$$

where $Q := [ q_1 \ \cdots \ q_n ] \in \mathbb{R}^{m \times n}$

If $\{ q_1, \cdots, q_n \}$ is a basis of $\mathbb{R}^m$, then $n = m$ and $r = 0$

i.e., $\quad v = \sum_{i=1}^{m} (q_i^T v) q_i = \sum_{i=1}^{n} (q_i^T q_i) v$

In fact, $v = QQ^T v$, i.e.,

$$QQ^T = I$$

Def. A square matrix $Q \in \mathbb{R}^{m \times m}$ is said to be orthogonal if

$$Q^T = Q^{-1}$$

$\uparrow$ <span style="color:green">should be called orthonormal</span>

i.e., $Q^T Q = QQ^T = I$

Note: If $Q = [q_1 \cdots q_n] \in \mathbb{R}^{m \times n}$ with $m > n$ and these vectors are orthonormal, then

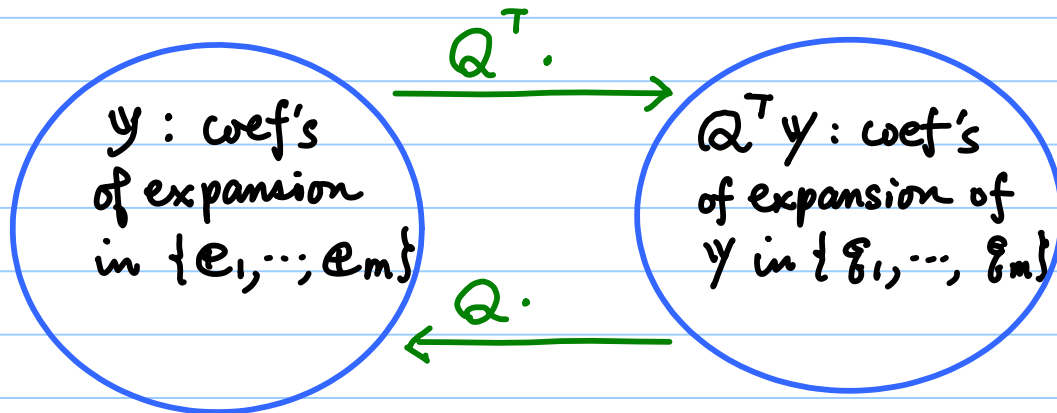it is always true that $Q^T Q = I_{n \times n}$ but $QQ^T \neq I_{m \times m}$ unless $m = n$

e.g.,

$$Q = \begin{bmatrix} 1/\sqrt{3} & 1/\sqrt{2} \\ 1/\sqrt{3} & 0 \\ 1/\sqrt{3} & -1/\sqrt{2} \end{bmatrix} \quad \text{then} \quad Q^T Q = \begin{bmatrix} 1 & 0 \\ 0 & 1 \end{bmatrix}$$

$$= I_{2 \times 2}$$

But $QQ^T = \begin{bmatrix} 1/\sqrt{3} & 1/\sqrt{2} \\ 1/\sqrt{3} & 0 \\ 1/\sqrt{3} & -1/\sqrt{2} \end{bmatrix} \begin{bmatrix} 1/\sqrt{3} & 1/\sqrt{3} & 1/\sqrt{3} \\ 1/\sqrt{2} & 0 & -1/\sqrt{2} \end{bmatrix}$

$$= \begin{bmatrix} 5/6 & 1/3 & -1/6 \\ 1/3 & 1/3 & 1/3 \\ -1/6 & 1/3 & 5/6 \end{bmatrix} \neq I_{3 \times 3}$$

why? $\Rightarrow$ Next lecture on Orthogonal Projector.

# ☆ Multiplication by an ortho. matrix



$$y : \text{coef's of expansion in } \{e_1, \cdots, e_m\} \xrightarrow{Q^T \cdot} Q^T y : \text{coef's of expansion of } y \text{ in } \{\xi_1, \cdots, \xi_m\}$$

$$\xleftarrow{Q \cdot}$$

Note that $\| y \| = \| Q^T y \|$ !

i.e., isometry!

why?

$$\| Q^T y \|^2 = (Q^T y)^T (Q^T y)$$
$$= y^T \underbrace{Q Q^T}_{I} y$$
$$= y^T y = \| y \|^2 \text{ !!}$$

Compare this with the general situation we discussed before: $A \in \mathbb{R}^{m \times n}$, nonsingular



$$y : \text{coef's of expansion in } \{e_1, \cdots, e_m\} \xrightarrow{A^{-1} \cdot} A^{-1} y : \text{coef's of expansion of } y \text{ in } \{a_1, \cdots, a_m\}$$

$$\xleftarrow{A \cdot}$$