



EECS 245 Fall 2025

Math for ML

lecture 15: Projections; Regression via Lin. Alg.
→ Read: Ch 2-10 (tons of new content),
Ch 3.1 (in progress)

Agenda

- Projecting onto span of multiple vectors
- The "normal equation"
- Simple linear regression using linear algebra

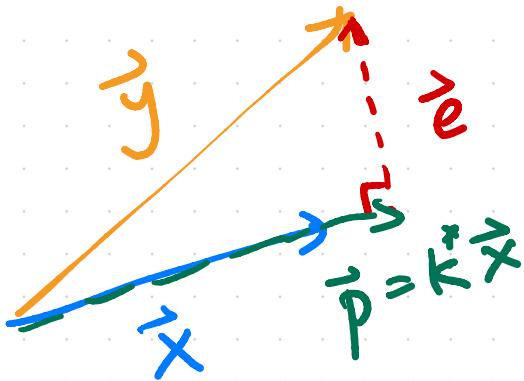
Announcements:

- ① Apply to be an IA! (dept form due Thurs, 245 form + video due Sat)
- ② HW 7 due Friday
 - Lab tomorrow: HW 7 work session

$$X_{n \times d} = \begin{bmatrix} | & | & | \\ \vec{x}^{(1)} & \vec{x}^{(2)} & \cdots & \vec{x}^{(d)} \\ | & | & | \end{bmatrix} \quad \vec{y} \in \mathbb{R}^n$$

Issue: \vec{y} is not (necessarily) in $\text{colsp}(X)$

Question: Among all vectors in $\text{colsp}(X)$,
which is "closest" to \vec{y} ?



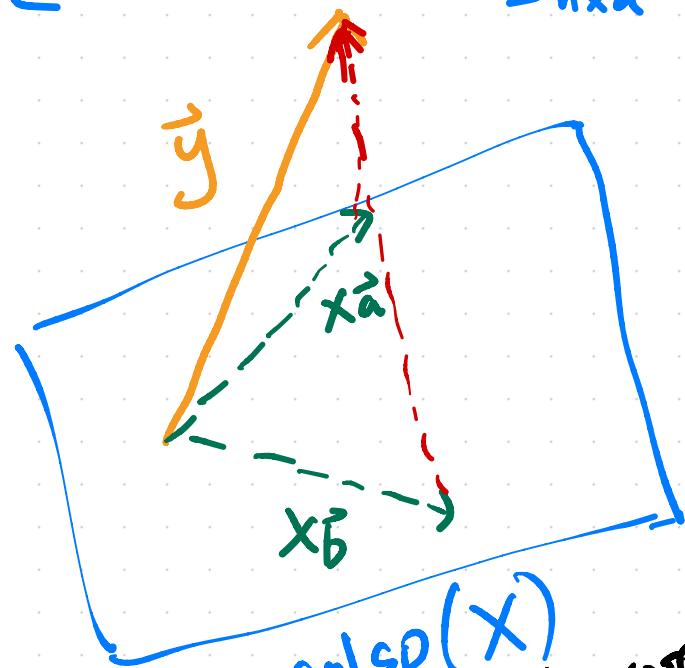
x^* chosen such that
 \vec{e} orthogonal to \vec{x} ;

that k^* minimized

$$f(k) = \|\vec{y} - k\vec{x}\|^2 = \|\vec{e}\|^2$$

$$\therefore k^* = \frac{\vec{x} \cdot \vec{y}}{\vec{x} \cdot \vec{x}}$$

$$X = \begin{bmatrix} | & | & | \\ x^{(1)} & x^{(2)} & \dots & x^{(d)} \\ | & | & | \end{bmatrix}_{n \times d}, \quad \vec{y} \in \mathbb{R}^n$$



$\text{colsp}(X)$
 = set of all lin comb's of X 's cols
 = set of all outputs of $X\vec{w}$,
 $\vec{w} \in \mathbb{R}^d$

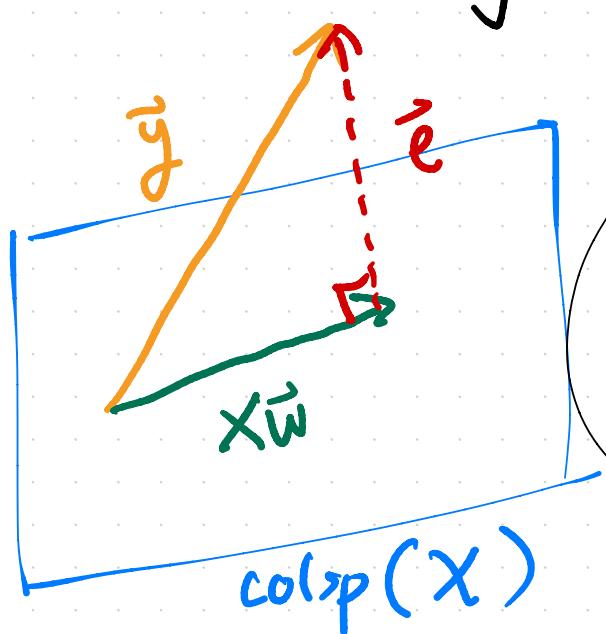
$$X w_d$$

$$X = \begin{bmatrix} | & | \\ 1 & 2 \\ 0 & -1 \\ 3 & -1 \end{bmatrix}_{3 \times 2}$$

key idea: to minimize

$$\|\vec{e}\|^2 = \|\vec{y} - X\vec{w}\|^2$$

pick \vec{w} so that $\vec{y} - X\vec{w}$ orthogonal to every vector in $\text{colsp}(X)$ function of \vec{w} only



← most important
diagram of
the semester

Well... which \vec{w} satisfies our goals?

Goal: Need $\vec{y} - X\vec{w}$ to be orthogonal to every col
of X

$$X = \begin{bmatrix} | & | & | \\ \vec{x}^{(1)} & \vec{x}^{(2)} & \dots & \vec{x}^{(d)} \\ | & | & | \end{bmatrix}_{n \times d}$$

$$\vec{x}^{(1)} \cdot (\vec{y} - X\vec{w}) = 0$$

$$\vec{x}^{(2)} \cdot (\vec{y} - X\vec{w}) = 0$$

⋮

$$X^T = \begin{bmatrix} \vec{x}^{(1)T} \\ \vec{x}^{(2)T} \\ \vdots \\ \vec{x}^{(d)T} \end{bmatrix} \quad d \times n$$

$$X^T(\vec{y} - X\vec{w}) = \begin{bmatrix} \vec{x}^{(1)} \cdot (\vec{y} - X\vec{w}) \\ \vec{x}^{(2)} \cdot (\vec{y} - X\vec{w}) \\ \vdots \\ \vec{x}^{(d)} \cdot (\vec{y} - X\vec{w}) \end{bmatrix} = \begin{bmatrix} 0 \\ 0 \\ \vdots \\ 0 \end{bmatrix} = \vec{0}_d$$

remember: we're looking for the "best" \vec{w}

\Rightarrow equivalent to finding \vec{w} that satisfies

$$X^T(\vec{y} - X\vec{w}) = \vec{0}$$

$$X^T\vec{y} - X^T X \vec{w} = \vec{0}$$

$$X^T X \vec{w} = X^T \vec{y}$$

system of d equations, d unknowns

$$\vec{w} = \begin{bmatrix} w_1 \\ w_2 \end{bmatrix}$$

the normal equations

$$X^T X$$

dot products
of all pairs
of cols

$$\underbrace{X^T X}_{d \times d} \vec{w} = X^T \vec{y}$$

Case (1) : $X^T X$ invertible (remember,
 which happens if and only if rank(X) = rank($X^T X$),
 X 's cols are linearly independent so $X^T X$ is invertible
 when rank($X^T X$)

Then, unique solution :

$$\vec{w}^* = (X^T X)^{-1} X^T \vec{y}$$

most important equation

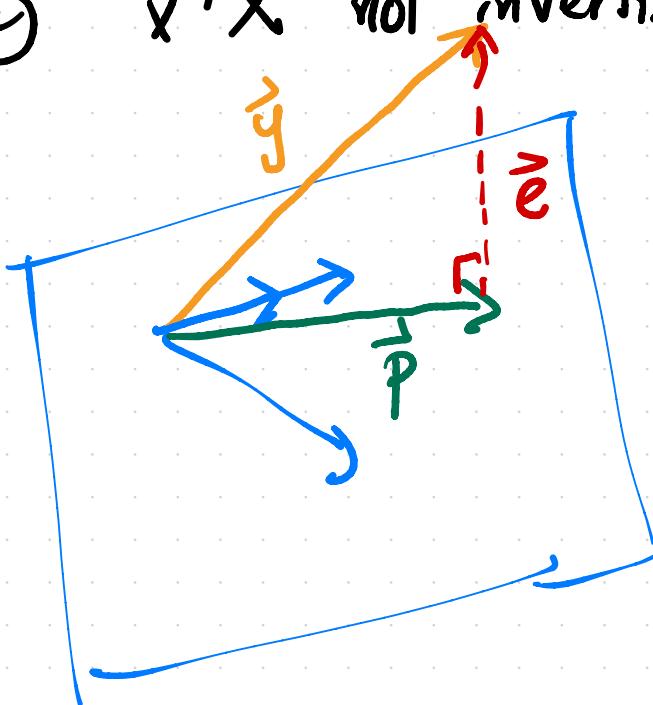
$$= \\ d \\ = \\ \text{rank}(X))$$

$$\chi^T X \vec{w} = X^T \vec{y}$$

Case ②

$X^T X$ not invertible

\rightarrow then, there are infinitely many \vec{w}^* that satisfy the normal equations, but all of them correspond to the same projection, \vec{p}



Example

$$X = \begin{bmatrix} 1 & 0 \\ 2 & 1 \\ 1 & -1 \\ 0 & -1 \end{bmatrix}, \quad \bar{y} = \begin{bmatrix} 1 \\ 0 \\ 4 \\ 5 \end{bmatrix}$$

Q: which lin comb of X 's cols is closest to \bar{y} ?

A: $X\vec{w}^*$, where

$$\vec{w}^* = (X^T X)^{-1} X^T \bar{y}$$

$$= \begin{bmatrix} 2 \\ -8/3 \end{bmatrix}$$

$$X^T X = \begin{bmatrix} 6 & 3 \\ 3 & 3 \end{bmatrix}$$

$$(X^T X)^{-1} = \frac{1}{9} \begin{bmatrix} 3 & -3 \\ -3 & 6 \end{bmatrix}$$

if $A = \begin{bmatrix} a & b \\ c & d \end{bmatrix}$,

then

$$A^{-1} = \frac{1}{ad - bc} \begin{bmatrix} d & -b \\ -c & a \end{bmatrix}$$

one question on the midterm
will be this exactly
→ practice it!

of all vectors of the form

$$w_1 \begin{bmatrix} 1 \\ 2 \\ -1 \\ 0 \end{bmatrix} + w_2 \begin{bmatrix} 0 \\ 1 \\ -1 \\ -1 \end{bmatrix}, \text{ the one that}$$

is closest to $\bar{y} = \begin{bmatrix} 1 \\ 6 \\ 4 \\ 5 \end{bmatrix}$ is

$$X = \begin{bmatrix} 1 & 0 \\ 2 & 1 \\ -1 & -1 \\ 0 & -1 \end{bmatrix} - 2 \begin{bmatrix} 1 \\ 2 \\ -1 \\ 0 \end{bmatrix} - \frac{8}{3} \begin{bmatrix} 0 \\ 1 \\ -1 \\ -1 \end{bmatrix}$$

orthogonality

$$\vec{p} = 2 \begin{bmatrix} 1 \\ 2 \\ -1 \\ 0 \end{bmatrix} - \frac{8}{3} \begin{bmatrix} 0 \\ 1 \\ 1 \\ -1 \end{bmatrix} = X \vec{w}^*$$

$\vec{e} = \vec{y} - \vec{p}$ is orthogonal to both cols of X

$$X = \begin{bmatrix} 1 & 2 \\ 3 & 2 \\ 0 & 2 \end{bmatrix}, Z = \begin{bmatrix} 1 & 0 \\ 0 & -1 \\ 0 & 1 \end{bmatrix}$$

$\vec{y} \in \mathbb{R}^3$

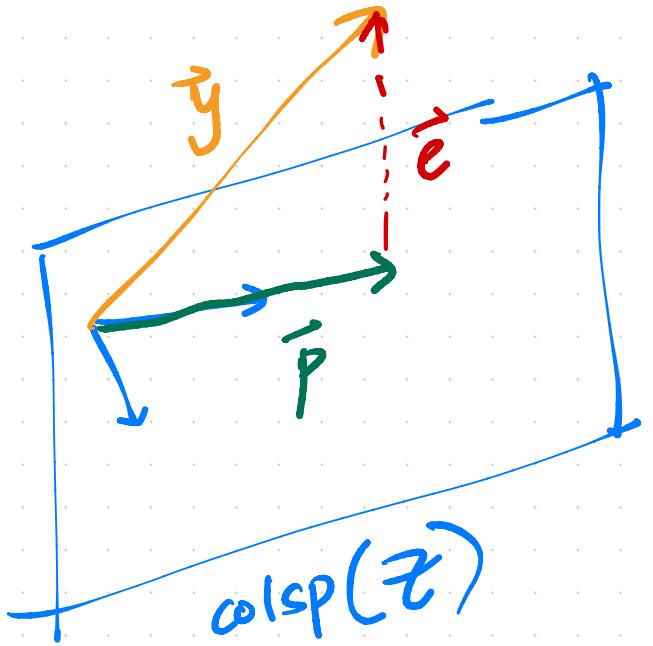
let \vec{P}_X be the orthogonal proj of \vec{y} onto $\text{colsp}(X)$

\vec{P}_Z $\cdots \cdots \cdots$ $\cdots \text{colsp}(Z)$

Why is it the case that the error vector for

\vec{P}_X sums to 0,

but not for \vec{P}_Z ?



$$X = \begin{bmatrix} 1 & 2 \\ 3 & 2 \\ 0 & 2 \end{bmatrix}$$

$$\vec{e}_x = \vec{y} - \vec{p}_x = \begin{bmatrix} e_1 \\ e_2 \\ e_3 \end{bmatrix}$$

Fact: \vec{e}_x orthogonal to all cols of X

$$\begin{bmatrix} e_1 \\ e_2 \\ e_3 \end{bmatrix} \cdot (\text{any lin comb of } X\text{'s cols}) = 0$$

$$\begin{bmatrix} e_1 \\ e_2 \\ e_3 \end{bmatrix} \cdot \begin{bmatrix} 2 \\ 2 \\ 2 \end{bmatrix} = 0 \rightarrow \cancel{2e_1 + 2e_2 + 2e_3 = 0}$$

$$e_1 + e_2 + e_3 = 0$$

orthogonal proj of \vec{y} onto $\text{colsp}(X)$

$$\vec{p} = X \vec{w}^* = \underbrace{X (X^T X)^{-1} X^T}_{P} \vec{y}$$

$$P = X (X^T X)^{-1} X^T$$

P "projection matrix" "hat matrix"

linear transformation of \vec{y}

read notes
Ch. 2.10

Linear regression!



3-step modeling process

① choose a model

$$h(x_i) = w_0 + w_1 x_i$$

② choose loss function

$$\text{Lsq}(y_i, h(x_i)) = (y_i - h(x_i))^2$$

③ minimize average loss



mean squared error

average squared loss

empirical risk

actual - pred

$$R_{\text{sq}}(w_0, w) = \frac{1}{n} \sum_{i=1}^n \left(y_i - \underbrace{(w_0 + w, x_i)}_{h(x_i)} \right)^2$$

this sum is
the (squared) norm
of error vector

definitions (Ch. 3.1)

given dataset

$$(x_1, y_1), (x_2, y_2), \dots, (x_n, y_n)$$

$$\vec{y} = \begin{bmatrix} y_1 \\ y_2 \\ \vdots \\ y_n \end{bmatrix}_n$$

"observation vector"

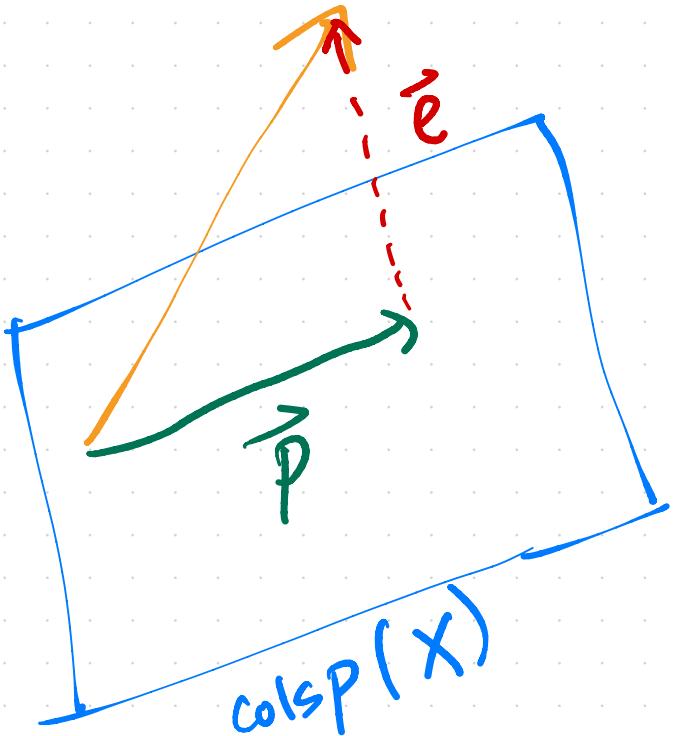
actual commute times

$$\vec{p} = \begin{bmatrix} w_0 + w_1 x_1 \\ w_0 + w_1 x_2 \\ \vdots \\ w_0 + w_1 x_n \end{bmatrix}_n = w_0 \begin{bmatrix} 1 \\ 1 \\ \vdots \\ 1 \end{bmatrix} + w_1 \begin{bmatrix} x_1 \\ x_2 \\ \vdots \\ x_n \end{bmatrix} = \begin{bmatrix} 1 \\ 1 \\ \vdots \\ 1 \end{bmatrix} \begin{bmatrix} x_1 \\ x_2 \\ \vdots \\ x_n \end{bmatrix} + \begin{bmatrix} w_0 \\ w_1 \end{bmatrix}$$

"prediction vector"

χ "design matrix"

"parameter vector" \vec{w}



$$\vec{e} = \vec{y} - \vec{p} = \begin{bmatrix} y_1 - (w_0 + w_1 x_1) \\ y_2 - (w_0 + w_1 x_2) \\ \vdots \\ y_n - (w_0 + w_1 x_n) \end{bmatrix}$$

Goal: minimize $\frac{1}{n} \|\vec{e}\|^2$

$$= \frac{1}{n} (y_1 - (w_0 + w_1 x_1))^2$$

$$+ \dots + (y_n - (w_0 + w_1 x_n))^2$$

$$= \frac{1}{n} \sum_{i=1}^n (y_i - (w_0 + w_1 x_i))^2$$

if you define

$$\vec{X} = \begin{bmatrix} 1 & x_1 \\ 1 & x_2 \\ \vdots & \vdots \\ 1 & x_n \end{bmatrix}$$

$$\vec{y} = \begin{bmatrix} y_1 \\ \vdots \\ y_n \end{bmatrix}$$

"design matrix"

then

$$\vec{\omega}^* = (\vec{X}^T \vec{X})^{-1} \vec{X}^T \vec{y} =$$

from before!

$$\begin{bmatrix} \vec{y} - r \frac{\sigma_y}{\sigma_x} \bar{x} \\ \frac{\sigma_y}{\sigma_x} \end{bmatrix}$$