

Lecture 3:
Least Squares

Given: vector of labels $\underline{y} \in \mathbb{R}^n$

matrix of features $X \in \mathbb{R}^{n \times p}$

Want: vector of weights $\underline{w} \in \mathbb{R}^p$

Assume: $n \geq p$

Rank(X) = p (X has p linearly independent columns)

$$\underbrace{\left. \begin{array}{c} \left. \begin{array}{c} \left. \begin{array}{c} X \\ \hline \end{array} \right\} p \\ \hline \end{array} \right\} n \\ \underline{w} \\ \hline \underline{y} \end{array} \right\} n =$$

\leftarrow n equations
 p unknowns

Let $\underline{x}_i \in \mathbb{R}^p$ be feature vector of i^{th} sample;
= (i^{th} row of X)^T

Let $\underline{X}_j \in \mathbb{R}^n$ be j^{th} feature for all samples
= j^{th} col of X

If $\hat{\underline{y}} = X\underline{w}$, then we have a system of n linear equations;

i^{th} equation: $y_i = w_1 x_{i1} + w_2 x_{i2} + \dots + w_p x_{ip}$

$$= \sum_{j=1}^p w_j x_{ij} = \langle \underline{w}, \underline{x}_i \rangle$$

\uparrow i^{th} row of X (transposed)

Side Note

(we might also consider $\hat{\underline{y}} = \underline{w}_0 + w_1 x_{01} + w_2 x_{02} + \dots + w_p x_{0p}$)

this can be done implicitly by letting $\underline{x}_0 = \begin{bmatrix} 1 \\ x_{01} \\ x_{02} \\ \vdots \\ x_{0p} \end{bmatrix}$, $\underline{w} = \begin{bmatrix} w_0 \\ w_1 \\ w_2 \\ \vdots \\ w_p \end{bmatrix} \in \mathbb{R}^{p+1}$

now our model is $\hat{\underline{y}} = \langle \underline{w}, \underline{x}_0 \rangle$, same as before!)

In general, $y \neq Xw$ for any w (because of modeling errors, noise)

Define residual $r_i = r_i(w) = y_i - \langle w, x_i \rangle = y_i - \hat{y}_i \Rightarrow r = [r_1, r_2, \dots, r_n]^T$

LEAST SQUARES ESTIMATION:

find w to minimize $\sum_{i=1}^n |r_i(w)|^2$

$$\langle r, r \rangle = \|r\|^2$$

Why least squares?

1) magnify effect of large errors

2) makes math easy

3) nice geometric interpretation

4) coincides with modeling $y = Xw + \underline{\varepsilon}$,
 $\underline{\varepsilon} = \text{Gaussian noise (later)}$

SPOILER ALERT!

best w (that minimizes sum of squared errors) is

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

and corresponding predicted labels are $\hat{y} = X\hat{w} = X(X^T X)^{-1} X^T y$

Today: when is this \hat{w} valid?

what is $(X^T X)^{-1}$?

where does this formula come from?

Span

The span of a set of vectors $\underline{x}_1, \underline{x}_2, \dots, \underline{x}_p \in \mathbb{R}^n$ is the set of vectors that can be written as a weighted sum of the \underline{x}_j 's :

$$\text{Span}(\underline{x}_1, \underline{x}_2, \dots, \underline{x}_p) = \left\{ \underline{y} \in \mathbb{R}^n : \underline{y} = \sum_{i=1}^p w_i \underline{x}_i \text{ for some } w_1, \dots, w_p \in \mathbb{R} \right\}$$

If $X = [\underline{x}_1, \underline{x}_2, \dots, \underline{x}_p]$, then $\text{range}(X) := \text{span}(\text{cols of } X) = \text{span}(\underline{x}_1, \dots, \underline{x}_p)$

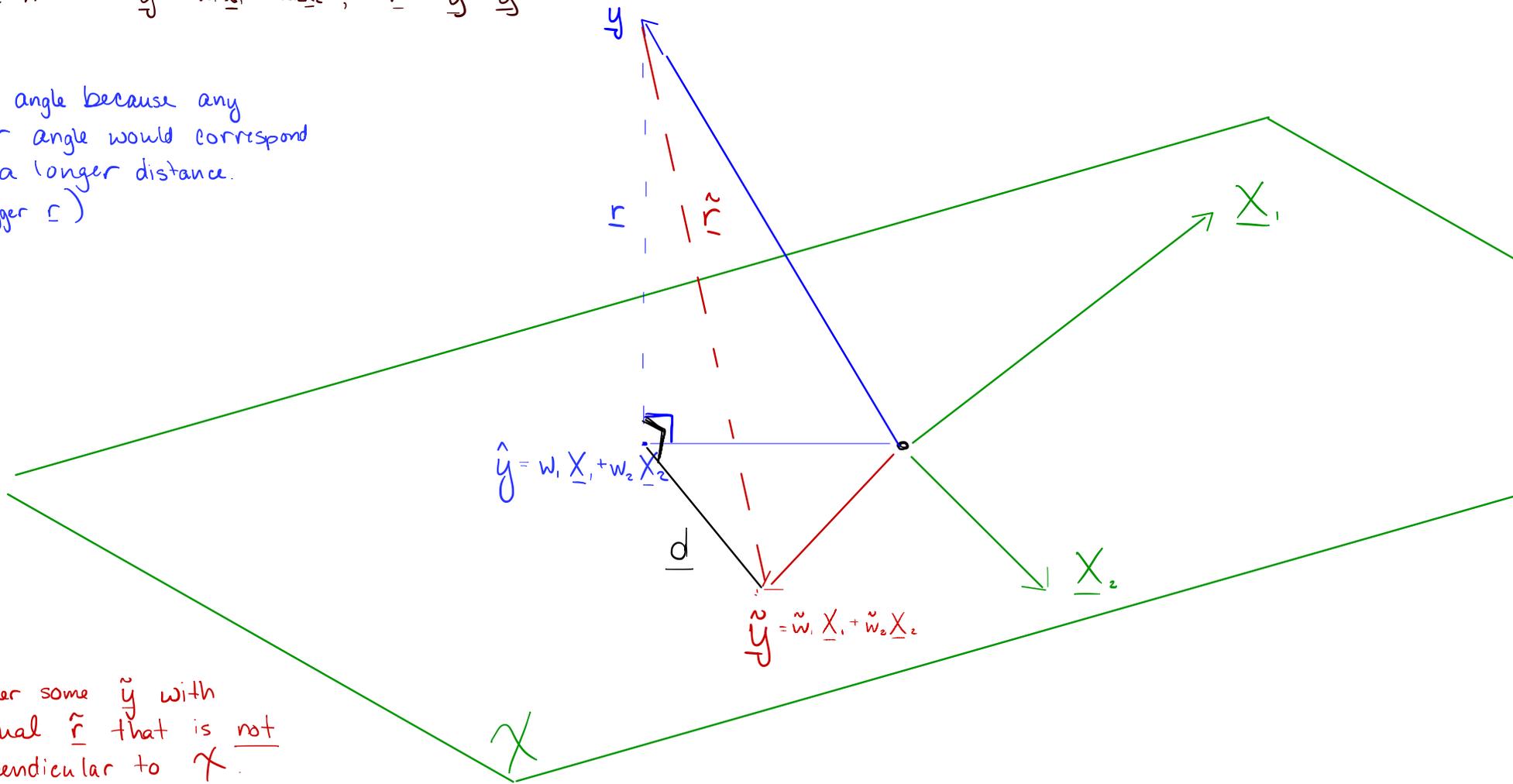
Ex. $\underline{x}_1 = \begin{bmatrix} 1 \\ 0 \\ 0 \end{bmatrix}$, $\underline{x}_2 = \begin{bmatrix} 0 \\ 1 \\ 0 \end{bmatrix}$. $\text{span}(\underline{x}_1, \underline{x}_2) =$ vectors of form $\begin{bmatrix} \alpha \\ \beta \\ 0 \end{bmatrix}$ for some α, β
i.e. vectors with zero in 3rd coordinate

Geometry of Least Squares

for general p , $\hat{y} = w_1 X_1 + w_2 X_2 + \dots + w_p X_p$

$p=2, n=3$. $\hat{y} = w_1 \underline{X}_1 + w_2 \underline{X}_2$, $r = y - \hat{y}$

Right angle because any other angle would correspond to a longer distance. (bigger ϵ)



Consider some \tilde{y} with residual \tilde{r} that is not perpendicular to X .

Then $\|\tilde{r}\|^2 = \|r\|^2 + \|d\|^2$ by Pythagorean Thm and so $\|\tilde{r}\|^2 > \|r\|^2$ so weights corresponding to \tilde{y} cannot be optimal (they don't make $\|\tilde{r}\|^2$ as small as possible.)

X = space of all vectors \hat{y} that can be written as $\hat{y} = \alpha \underline{X}_1 + \beta \underline{X}_2$ for some $\alpha, \beta \in \mathbb{R}$. called "span" of cols of X . y may not be in this space

$\hat{\underline{w}}$ = "argument \underline{w} that minimizes" $\sum_{i=1}^n r_i^2(\underline{w})$

$$= \operatorname{argmin}_{\underline{w}} \sum_{i=1}^n r_i^2(\underline{w}) = \operatorname{argmin}_{\underline{w}} \langle r, r \rangle = \operatorname{argmin}_{\underline{w}} \|r\|_2^2$$

$$\text{Let } \underline{r} := r(\hat{\underline{w}}) = \begin{bmatrix} r_1(\hat{\underline{w}}) \\ \vdots \\ r_n(\hat{\underline{w}}) \end{bmatrix}$$

We know that $\underline{r} = \underline{y} - X\hat{\underline{w}}$ is perpendicular/orthogonal to span of columns of X

This implies $\underline{X}_i^T \underline{r} = 0$ for each column i of X

$$\Rightarrow X^T \underline{r} = \underline{0} \leftarrow \text{vector of zeros}$$

$$\Rightarrow X^T (\underline{y} - X\hat{\underline{w}}) = \underline{0}$$

$$\Rightarrow \hat{\underline{w}} \text{ is solution to linear system of equations } X^T \underline{y} = X^T X \underline{w}$$

Two vectors $\underline{u}, \underline{v}$ are
orthogonal if $\langle \underline{u}, \underline{v} \rangle = 0$

WHAT CAN WE SAY ABOUT \underline{w} SATISFYING $X^T \underline{y} = X^T X \underline{w}$?

- does it exist ?
- is it unique ?

↳ "yes" to both when columns of X are linearly independent

Consider following linear systems For each, how many solutions are there?

(zero, one, or many) If one or more solutions exist, find one or more Why do different cases have different numbers of solutions?

a) $3w_1 + 2w_2 = 1$
 $3w_1 + w_2 = 0$

$\Rightarrow w_2 = -3w_1$
 $-w_2 + 2w_2 = w_2 = 1$
 $\Rightarrow w_1 = -\frac{1}{3} \Rightarrow \text{one soln}$

$\text{rank}(X) = 2$
 \downarrow
 $\underbrace{\begin{bmatrix} 3 & 2 \\ 3 & 1 \end{bmatrix}}_X \underbrace{\begin{bmatrix} w_1 \\ w_2 \end{bmatrix}}_w = \underbrace{\begin{bmatrix} 1 \\ 0 \end{bmatrix}}_y$

b) $\begin{bmatrix} 3 & 2 \\ 3 & 1 \\ 3 & 3 \end{bmatrix} \begin{bmatrix} x_1 \\ x_2 \end{bmatrix} = \begin{bmatrix} 1 \\ 0 \\ 2 \end{bmatrix} \Rightarrow 1 \text{ soln: } x_1 = -\frac{1}{3}, x_2 = 1$

c) $\left. \begin{array}{l} 3x_1 + 2x_2 = 1 \\ 3x_1 + x_2 = 0 \\ 2x_1 + 2x_2 = 2 \end{array} \right\} \begin{array}{l} x_1 = -\frac{1}{3}, x_2 = 1 \\ 2(-\frac{1}{3}) + 2(1) \neq 2 \end{array} \Rightarrow 0 \text{ solutions}$

$\begin{bmatrix} 3 & 2 \\ 3 & 1 \\ 2 & 2 \end{bmatrix} \begin{bmatrix} x_1 \\ x_2 \end{bmatrix} = \begin{bmatrix} 1 \\ 0 \\ 2 \end{bmatrix}$

d) $3x_1 + x_2 = 0 \Rightarrow \infty \text{ solutions}$

e) $\begin{array}{l} 3x_1 + x_2 = 1 \\ 6x_1 + 2x_2 = 2 \end{array} \Leftrightarrow \begin{bmatrix} 3 & 1 \\ 6 & 2 \end{bmatrix} \begin{bmatrix} x_1 \\ x_2 \end{bmatrix} = \begin{bmatrix} 1 \\ 2 \end{bmatrix}$
 $\Rightarrow \infty \text{ solns}$ (rank-1 matrix!)

if X is $n \times p$, $p \leq n$, and $\text{rank}(X) < p$,
 then $X^T X$ not invertible AND no unique solution

Linear Independence

A collection of vectors $\underline{v}_1, \underline{v}_2, \dots, \underline{v}_p \in \mathbb{R}^n$ is linearly independent when

$$\sum_{i=1}^p \alpha_i \underline{v}_i = \underline{0} \quad \text{if and only if} \quad \alpha_i = 0 \quad \text{for } i=1, 2, \dots, p$$

That is, any weighted sum of the vectors is nonzero unless all the weights are zero

EX $n=3$ $p=2$ $\underline{v}_1 = \begin{bmatrix} 1 \\ 0 \\ 0 \end{bmatrix}$, $\underline{v}_2 = \begin{bmatrix} 0 \\ 1 \\ 1 \end{bmatrix}$ \Rightarrow yes, linearly independent (LI)

note $\alpha_1 \underline{v}_1 + \alpha_2 \underline{v}_2 = \begin{bmatrix} \alpha_1 \\ \alpha_2 \\ \alpha_2 \end{bmatrix}$ This can only be the zero vector if $\alpha_1 = \alpha_2 = 0$

EX $n=3$ $p=3$ $\underline{v}_1 = \begin{bmatrix} 1 \\ 0 \\ 0 \end{bmatrix}$, $\underline{v}_2 = \begin{bmatrix} 0 \\ 1 \\ 1 \end{bmatrix}$, $\underline{v}_3 = \begin{bmatrix} 0 \\ 1 \\ 0 \end{bmatrix}$ \Rightarrow yes, linearly independent (LI)

$\alpha_1 \underline{v}_1 + \alpha_2 \underline{v}_2 + \alpha_3 \underline{v}_3 = \begin{bmatrix} \alpha_1 \\ \alpha_2 + \alpha_3 \\ \alpha_2 \end{bmatrix}$ this = 0 only if $\alpha_1 = \alpha_2 = \alpha_3 = 0$

$$\text{Ex } n=3 \\ p=4 \quad \underline{v}_1 = \begin{bmatrix} 1 \\ 0 \\ 0 \end{bmatrix}, \quad \underline{v}_2 = \begin{bmatrix} 0 \\ 1 \\ 1 \end{bmatrix}, \quad \underline{v}_3 = \begin{bmatrix} 0 \\ 1 \\ 0 \end{bmatrix}, \quad \underline{v}_4 = \begin{bmatrix} 1 \\ 0 \\ 1 \end{bmatrix}$$

note $\underline{v}_4 = \underline{v}_1 + \underline{v}_2 - \underline{v}_3$ (we can write one vector as linear combination (weighted sum) of others This implies linear dependence)

$$\alpha_1 \underline{v}_1 + \alpha_2 \underline{v}_2 + \alpha_3 \underline{v}_3 + \alpha_4 \underline{v}_4 = \begin{bmatrix} \alpha_1 + \alpha_4 \\ \alpha_2 + \alpha_3 \\ \alpha_2 + \alpha_4 \end{bmatrix} \Rightarrow \text{if } \alpha_1 = -\alpha_4 = \alpha_2 = -\alpha_3, \text{ then}$$

$$\alpha_1 \underline{v}_1 + \alpha_2 \underline{v}_2 + \alpha_3 \underline{v}_3 + \alpha_4 \underline{v}_4 = 0$$

$$\Rightarrow \text{NOT linearly independent}$$

Linear independence $\Rightarrow n \geq p$

$p > n \Rightarrow$ Linear dependence

Matrix rank number of linearly independent columns = # linearly independent rows

if $X^T = [\underline{x}_1 \quad \underline{x}_2 \quad \dots \quad \underline{x}_n] \in \mathbb{R}^{p \times n}$, then $\text{rank}(X) \leq \min(p, n)$

$$\text{ex } X = \begin{bmatrix} 1 & 0 \\ 0 & 1 \\ 0 & 1 \end{bmatrix} \Rightarrow \text{rank}(X) = 2$$

$$\text{ex } X = \begin{bmatrix} 1 & 0 & 1 \\ 0 & 1 & 1 \\ 0 & 1 & 1 \end{bmatrix} \Rightarrow \text{rank}(X) = 2$$

Matrix Inverse

for a square matrix A , its inverse A^{-1} is a square matrix that satisfies:

$$AA^{-1} = A^{-1}A = I = \begin{bmatrix} 1 & & & \\ & \ddots & & \\ & & \ddots & \\ & & & 1 \end{bmatrix}$$

$$\text{Ex. } A = \begin{bmatrix} 1/4 & 0 \\ 0 & 2 \end{bmatrix} \Rightarrow A^{-1} = \begin{bmatrix} 4 & 0 \\ 0 & 1/2 \end{bmatrix}$$

$$\text{Ex. } A = \begin{bmatrix} 1 & 2 \\ 3 & 4 \end{bmatrix} \Rightarrow A^{-1} = \begin{bmatrix} -2 & 1 \\ 3/2 & -1/2 \end{bmatrix}$$

Not all matrices have inverses.

Specifically, A only has an inverse if it is full rank

$$\text{Ex: } A = \begin{bmatrix} 1 & 0 \\ 0 & 0 \end{bmatrix}, \quad A = \begin{bmatrix} 1 & 2 \\ 2 & 4 \end{bmatrix} \quad \text{have no inverse.}$$

$$X \in \mathbb{R}^{n \times p}, \quad n \geq p, \quad \text{rank}(X) = p \Rightarrow \text{rank}(X^T X) = p \Rightarrow X^T X \text{ has inverse}$$

Recall from earlier:

$\hat{\underline{w}}$ = "argument \underline{w} that minimizes" $\sum_{i=1}^n r_i^2(\underline{w}) = \operatorname{argmin}_{\underline{w}} \sum_{i=1}^n r_i^2(\underline{w})$

$$\text{Let } \hat{\underline{r}} := r(\hat{\underline{w}}) = \begin{bmatrix} r_1(\hat{\underline{w}}) \\ \vdots \\ r_n(\hat{\underline{w}}) \end{bmatrix}$$

We know that $\hat{\underline{r}} = \underline{y} - X\hat{\underline{w}}$ is perpendicular/orthogonal to span of columns of X

This implies $\underline{X}_i^T \hat{\underline{r}} = 0$ for each column i of X

$$\Rightarrow X^T \hat{\underline{r}} = \underline{0} \quad \leftarrow \text{vector of zeros}$$

$$\Rightarrow X^T (\underline{y} - X\hat{\underline{w}}) = \underline{0}$$

$$\Rightarrow \hat{\underline{w}} \text{ is solution to linear system of equations } X^T \underline{y} = X^T X \underline{w}$$

Two vectors $\underline{u}, \underline{v}$ are
orthogonal if $\langle \underline{u}, \underline{v} \rangle = 0$

So if $X^T X$ is invertible (ie. if $X^T X$ is full-rank, which occurs if $\operatorname{rank}(X) = p \leq n$) then there is a unique solution:

$$\hat{\underline{w}} = (X^T X)^{-1} X^T \underline{y}$$

$$\begin{aligned} \Rightarrow \hat{\underline{y}} &= X \hat{\underline{w}} \\ &= X (X^T X)^{-1} X^T \underline{y} \end{aligned}$$

$P_X := X (X^T X)^{-1} X^T$ is called a **Projection Matrix**
because $P_X \underline{y}$ projects \underline{y} onto $\operatorname{range}(X)$

$$\Sigma_X \quad X = \begin{bmatrix} u_1 \\ u_2 \\ u_3 \end{bmatrix} \begin{bmatrix} v_1 & v_2 & v_3 \end{bmatrix} = \begin{bmatrix} u_1 v_1 & u_1 v_2 & u_1 v_3 \\ u_2 v_1 & u_2 v_2 & u_2 v_3 \\ u_3 v_1 & u_3 v_2 & u_3 v_3 \end{bmatrix} \rightarrow \text{Rank}(X) = 1$$

Recall the outer product representation of matrix product

$$UV = \begin{bmatrix} | & | & | \\ U_1 & U_2 & U_r \\ | & | & | \end{bmatrix} \begin{bmatrix} - V_1^T - \\ - V_2^T - \\ \vdots \\ - V_r^T - \end{bmatrix}$$

$$= \begin{array}{c} | \\ U_1 \\ | \end{array} \begin{array}{c} \text{---} \\ V_1^T \\ \text{---} \end{array} + \begin{array}{c} | \\ U_2 \\ | \end{array} \begin{array}{c} \text{---} \\ V_2^T \\ \text{---} \end{array} + \dots + \begin{array}{c} | \\ U_r \\ | \end{array} \begin{array}{c} \text{---} \\ V_r^T \\ \text{---} \end{array}$$

$$= \sum_{k=1}^r U_k V_k^T = \text{sum of rank-1 matrices} \Rightarrow \text{rank}(UV) = r \text{ if } U_k \text{'s are LI} \\ \text{and } V_k \text{'s are LI}$$