

Warmup: Suppose  $f(w) = 3w^2 - 5w + 2$

what value of  $w$  minimizes  $f$ ?

$$\frac{df}{dw} = 6w - 5 = 0$$

$$w = 5/6$$

## Lecture 5: Least Squares

given:

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

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

want:

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

assume:

$n \geq p$ ,  $\text{Rank}(X) = p$  (have  $p$  linearly independent columns)

if  $y = X\underline{w}$ , then have system of  $n$  linear equations

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

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

$\curvearrowleft i^{\text{th}}$  row of  $X$

(2)

In general,  $\underline{y} \neq X\underline{w}$  for any  $\underline{w}$

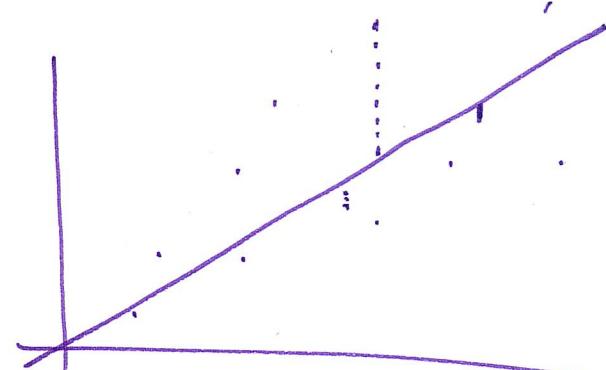
(because modeling errors, noise)

Define residual  $r_i = y_i - \langle \underline{w}, \underline{x}_{\cdot i} \rangle$

Goal: find  $\underline{w}$  to minimize  $\sum_{i=1}^n |r_i|^2$  (sum of squared residuals/errors)

Why sum of squared errors?

- a) magnify effect of large errors
- b) make math easy (can compute derivatives)
- c) nice geometric interpretation
- d) coincides w/ modeling  $\underline{y} = X\underline{w} + \underline{\epsilon}$ ,  $\underline{\epsilon}$  = Gaussian noise  
(not this course)



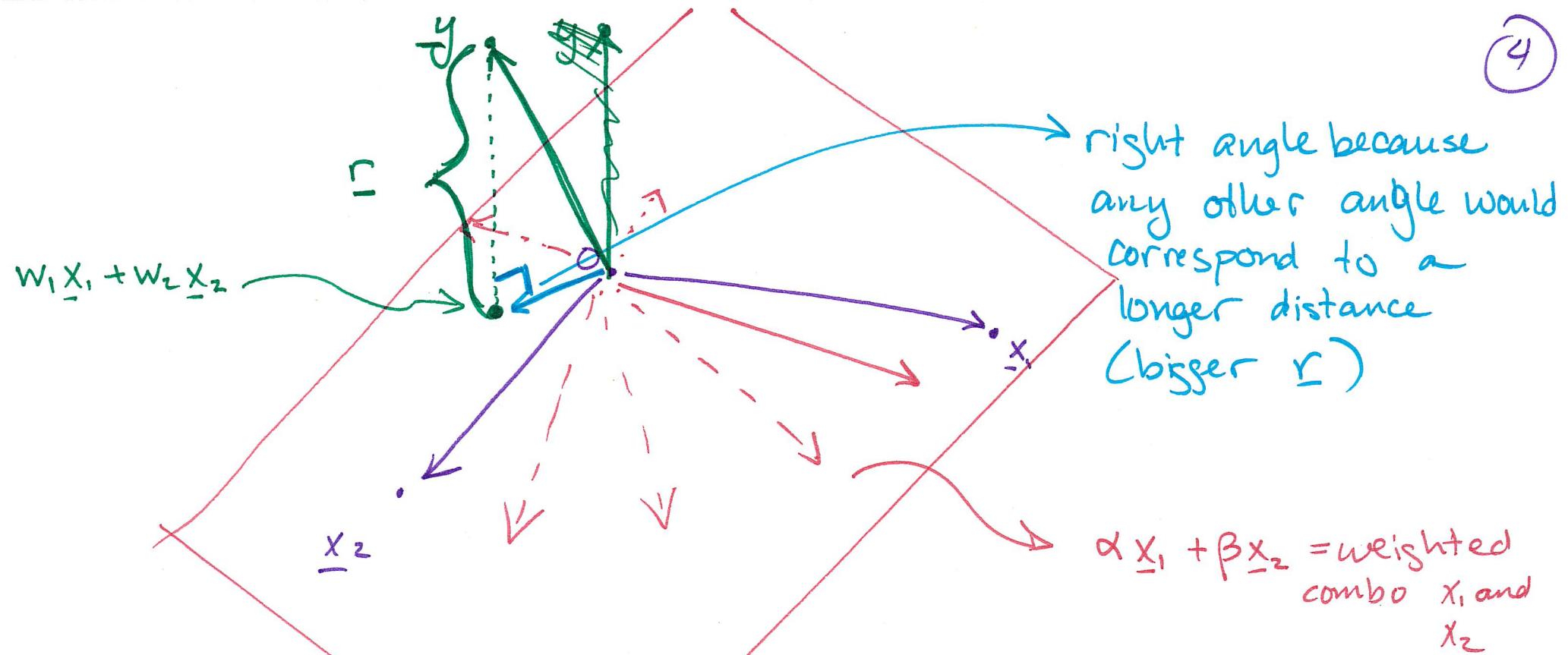
First, Geometry.

$$p = 2, n = 3$$

$$\begin{bmatrix} y_1 \\ y_2 \\ y_3 \end{bmatrix} \approx \begin{bmatrix} x_{11} \\ x_{21} \\ x_{31} \end{bmatrix} w_1 + \begin{bmatrix} x_{12} \\ x_{22} \\ x_{32} \end{bmatrix} w_2 , \quad r = \underline{y} - \underline{w}_1 \underline{x}_1 - \underline{w}_2 \underline{x}_2$$

$\uparrow$  1<sup>st</sup> col of  $X$

(4)



What is the point in  $\mathcal{X}$  that has the shortest distance to  $y$ ?

$\mathcal{X}$

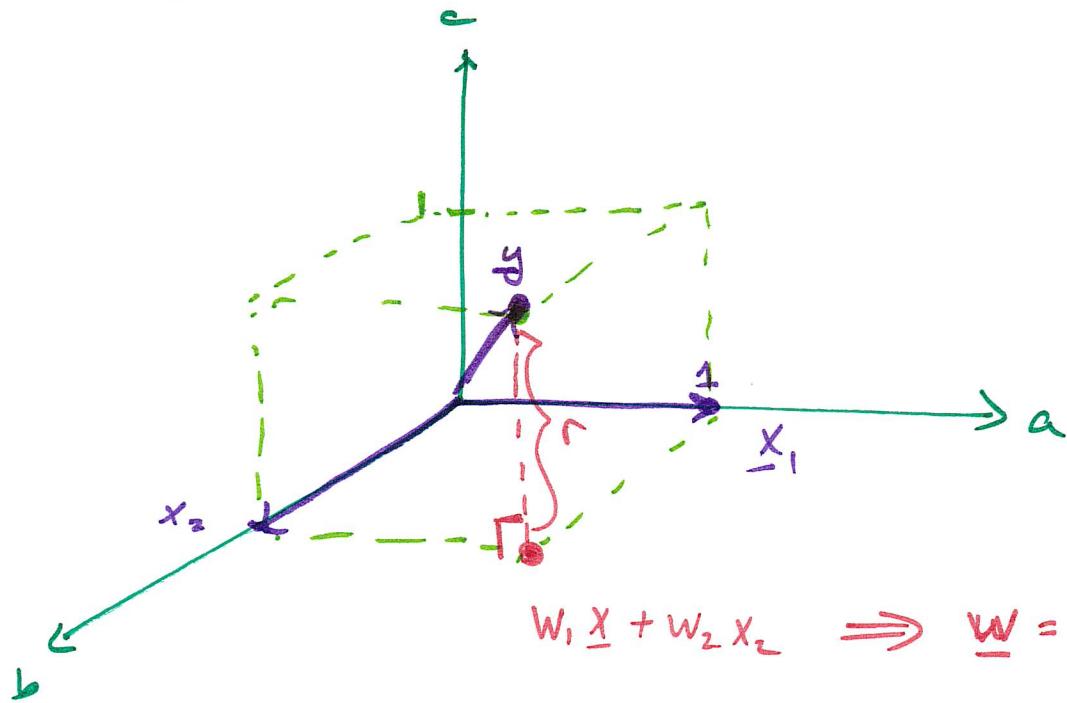
$\mathcal{X} = \text{space of all vectors that can be written as } \alpha \underline{x}_1 + \beta \underline{x}_2 \text{ for some } \alpha, \beta \in \mathbb{R}$

$= \text{span of cols of } X$

$(y \text{ may not lie in this space})$

(5)

$$ex) \quad \underline{x}_1 = \begin{bmatrix} 1 \\ 0 \\ 0 \end{bmatrix}, \quad \underline{x}_2 = \begin{bmatrix} 0 \\ 1 \\ 0 \end{bmatrix}, \quad y = \begin{bmatrix} 1 \\ 1 \end{bmatrix}$$



$$w_1 \underline{x}_1 + w_2 \underline{x}_2 \Rightarrow \underline{w} = \begin{bmatrix} 1 \\ -1 \\ 0 \end{bmatrix}, \quad \underline{X_w} = \begin{bmatrix} 1 \\ -1 \\ 0 \end{bmatrix}$$

(6)

Algebraic approach:

goal: choose  $\underline{w}$  to minimize  $\sum_{i=1}^n r_i^2$

$$\hat{\underline{w}} = \arg \min_{\underline{w}} \sum_{i=1}^n r_i^2$$

↑ argument that minimizes objective

$$= \arg \min_{\underline{w}} \sum_{i=1}^n \left( y_i - \sum_{j=1}^p w_j x_{ij} \right)^2$$

$$= \arg \min_{\underline{w}} \|\underline{r}\|_2^2$$

$$= \arg \min_{\underline{w}} \underbrace{\|y - X\underline{w}\|_2^2}_{f(\underline{w})}$$

$$\left( \sum_i a_i^2 \right)^{1/2} = \|\underline{a}\|_2$$

Approach 1: use geometry.

know  $\hat{r} = \underline{y} - \underline{X}\hat{\underline{w}}$  is orthogonal / perpendicular to plane / span of cols of  $\underline{X}$ .

$\leftarrow$   $i^{\text{th}}$  col of  $\underline{X}$

$$\underline{x}_i^T \cancel{\hat{r}} = 0 \quad \text{for each } i$$

$$\underline{X}^T \hat{r} = 0$$

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

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

{ two vectors  $u, v$  are orthogonal if  $\langle u, v \rangle = u^T v = 0$

$$\hat{r} \in \mathbb{R}^n$$

$$\underline{x}_i \in \mathbb{R}^n \quad (\text{col of } \underline{X})$$

$$\underline{x}_i \in \mathbb{R}^p \quad (\text{row of } \underline{X})$$

Matrix inverse:

for a square matrix  $A$ , ~~if~~ its inverse  $A^{-1}$  is a square matrix satisfies

$$AA^{-1} = A^{-1}A = I \Rightarrow \begin{bmatrix} 1 & & 0 \\ & \ddots & \\ 0 & & 1 \end{bmatrix}$$

e.g.  $A = \begin{bmatrix} \frac{1}{4} & 0 \\ 0 & 2 \end{bmatrix} \Rightarrow A^{-1} = \begin{bmatrix} 4 & 0 \\ 0 & \frac{1}{2} \end{bmatrix}$

$$A = \begin{bmatrix} 1 & 2 \\ 3 & 4 \end{bmatrix} \Rightarrow A^{-1} = \begin{bmatrix} -2 & 1 \\ \frac{3}{2} & -\frac{1}{2} \end{bmatrix}$$

not all matrices have inverses:

$$A = \begin{bmatrix} 1 & 2 \\ 2 & 4 \end{bmatrix} \Rightarrow \text{no inverse}$$

Back to Least squares:

$$\hat{\underline{w}} \text{ satisfies } \underline{X}^T \underline{y} = \underline{X}^T \underline{X} \hat{\underline{w}}$$

so if  $\underline{X}^T \underline{X}$  is invertible ( $(\underline{X}^T \underline{X})^{-1}$  exists)

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

Approach 2: use vector calculus

if we have a vector  $\underline{w}$  and a function of it  $f(\underline{w})$

then the gradient (vector version of derivative) is

$$\nabla_{\underline{w}} f(\underline{w}) = \begin{bmatrix} df(\underline{w})/dw_1 \\ df(\underline{w})/dw_2 \\ \vdots \\ df(\underline{w})/dw_p \end{bmatrix} = \text{vector same size as } \underline{w}$$

ex.  $f(\underline{w}) = \langle \underline{a}, \underline{w} \rangle = \underline{w}^T \underline{a} = \sum_{i=1}^p w_i a_i$

$$\nabla_{\underline{w}} \underline{w}^T \underline{a} = \begin{bmatrix} a_1 \\ a_2 \\ \vdots \\ a_p \end{bmatrix} = \underline{a}$$

(analog  $f(w) = aw$ ,  
 $\frac{df}{dw} = a$ )

$$\nabla_{\underline{w}} \underline{w}^T A \underline{w} = A \underline{w} + A^T \underline{w} ; \text{ if } A \text{ is symmetric}$$

$\uparrow$   
d x p matrix

$(A = A^T)$ , then

$$\nabla_{\underline{w}} \underline{w}^T A \underline{w} = 2A \underline{w}$$

(analog  $f(w) = aw^2$ ,  
 $\frac{df}{dw} = 2aw$ )

Least squares:

$$\hat{\underline{w}} = \arg \min_{\underline{w}} \| \underline{y} - \underline{X}\underline{w} \|_2^2 \quad \longrightarrow \| \underline{a} \|_2^2 = \underline{a}^T \underline{a}$$

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

$$= \underbrace{\underline{y}^T \underline{y} - \underline{y}^T \underline{X} \underline{w}}_{= \underline{w}^T \underline{X}^T \underline{y}} - \underline{w}^T \underline{X}^T \underline{y} + \underline{w}^T \underline{X}^T \underline{X} \underline{w}$$

$$= \underline{y}^T \underline{y} - 2\underbrace{\underline{w}^T \underline{X}^T \underline{y}}_a + \underbrace{\underline{w}^T \underline{X}^T \underline{X} \underline{w}}_A$$

$$\nabla_{\underline{w}} f(\underline{w}) = 0 - 2\underbrace{\underline{X}^T \underline{y}}_a + 2\underbrace{\underline{X}^T \underline{X} \underline{w}}_A = 0$$

$\hat{\underline{w}}$  solves  $\underline{X}^T \underline{y} = \underline{X}^T \underline{X} \underline{w}$

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

$\underline{X}^T \underline{X}$  is  
symmetric  
for any  $\underline{X}$

Does  $(X^T X)^{-1}$  exist?

Thm 3.10: If columns of  $X$  are linearly independent, then  $X^T X$  is invertible / full rank / non-singular and there exists a least squares solution

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

Proof:

$X^T X$  is symmetric, so  $X^T X$  is invertible if it is positive definite — i.e. for all  $\underline{z} \neq 0$ ,

$$\underline{z}^T X^T X \underline{z} > 0$$

Let  $\tilde{\underline{z}} = X \underline{z}$ . Then  $\underbrace{\tilde{\underline{z}}^T}_{\tilde{\underline{z}}^T} \underbrace{X^T X \tilde{\underline{z}}}_{\tilde{\underline{z}}} = \tilde{\underline{z}}^T \tilde{\underline{z}} = \|\tilde{\underline{z}}\|_2^2$

Recall:

$\|\tilde{z}\|_2 \geq 0$  for all  $\tilde{z}$ , and  $\|\tilde{z}\|_2 = 0$  if and only if  $\tilde{z} = 0$

now  $\tilde{z} = X z$ . can  $\tilde{z} = 0$  when  $z \neq 0$  ?

no. because  $\tilde{z}$  = weighted sum of cols of  $X$

we assumed cols of  $X$  were linearly independent

therefore  $\tilde{z} = 0$  implies  $z = 0$

thus

$$z^\top X^\top X z > 0 \quad \text{for all } z \neq 0.$$

$\Rightarrow X^\top X$  is invertible.

Is  $\hat{w}$  the Least squares solution?

If it is, then for any (other)  $w$

$$\|y - X\hat{w}\|_2^2 \leq \|y - Xw\|_2^2 \iff \|\hat{r}\|_2^2 \leq \|r\|_2^2$$

Let  $\hat{r} = y - X\hat{w}$ ,  $r = y - Xw$

$$r = y - Xw = \underbrace{y - X\hat{w}}_{\hat{r}} + X\hat{w} - Xw$$

$$r = \hat{r} + X(\hat{w} - w)$$

$$\|r\|_2^2 = r^T r = (\hat{r} + X(\hat{w} - w))^T (\hat{r} + X(\hat{w} - w))$$

$$\begin{aligned} &= \hat{r}^T \hat{r} + \cancel{\hat{r}^T X(\hat{w} - w)} + (\hat{w} - w)^T \cancel{X^T \hat{r}} + \underbrace{(\hat{w} - w)^T X^T X (\hat{w} - w)}_{\geq 0 \text{ b/c } X^T X \text{ p.d.}} \\ &\geq \|\hat{r}\|_2^2 + 0 + 0 \end{aligned}$$

$$\|r\|_2^2 \geq \|\hat{r}\|_2^2, \text{ and only equal if } \hat{w} = w \Rightarrow \hat{w} = (X^T X)^{-1} X^T y$$