

Lecture 3: Least squares and geometry

Mathematical Foundations of Machine Learning
University of Chicago

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: $n \geq p$

$\text{Rank}(X) = p$ (X has p

linearly independent columns)

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

$$\begin{aligned} i^{\text{th}} \text{ equation: } y_i &= w_1 x_{i1} + w_2 x_{i2} + \dots + w_p x_{ip} \\ &= \sum_{j=1}^p w_j x_{ij} = \langle \underline{w}, \underline{x}_i \rangle \end{aligned}$$

↑ i^{th} row of X (transposed)

In general, $y \neq X\underline{w}$ for any \underline{w} (because of modeling errors, noise)

Define residual $r_i = r_i(\underline{w}) = y_i - \langle \underline{w}, \underline{x}_i \rangle$

LEAST SQUARES ESTIMATION:

$$\begin{aligned} \text{find } \underline{w} \text{ to minimize } &\sum_{i=1}^n |r_i(\underline{w})|^2 \\ &= \|\underline{r}\|^2 \end{aligned}$$

$$n \left\{ \begin{array}{c} X \\ \vdots \\ \underline{w} \end{array} \right. = \left\{ \begin{array}{c} y \\ \vdots \\ \underline{y} \end{array} \right. \quad \left. \begin{array}{c} n \\ \vdots \\ p \end{array} \right\} \quad \left. \begin{array}{c} \leftarrow \\ n \text{ equations} \\ p \text{ unknowns} \end{array} \right.$$

let $x_i \in \mathbb{R}^p$ be feature vector of i^{th} sample;
 $= (i^{\text{th}} \text{ row of } X)^T$

let $x_j \in \mathbb{R}^n$ be j^{th} feature for all samples
 $= j^{\text{th}}$ col of X

Why least squares?

- 1) magnify effect of large errors
- 2) makes math easy
- 3) nice geometric interpretation
- 4) coincides with modeling $y = X\underline{w} + \underline{\epsilon}$,
 $\underline{\epsilon}$ = Gaussian noise (later)

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 :

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

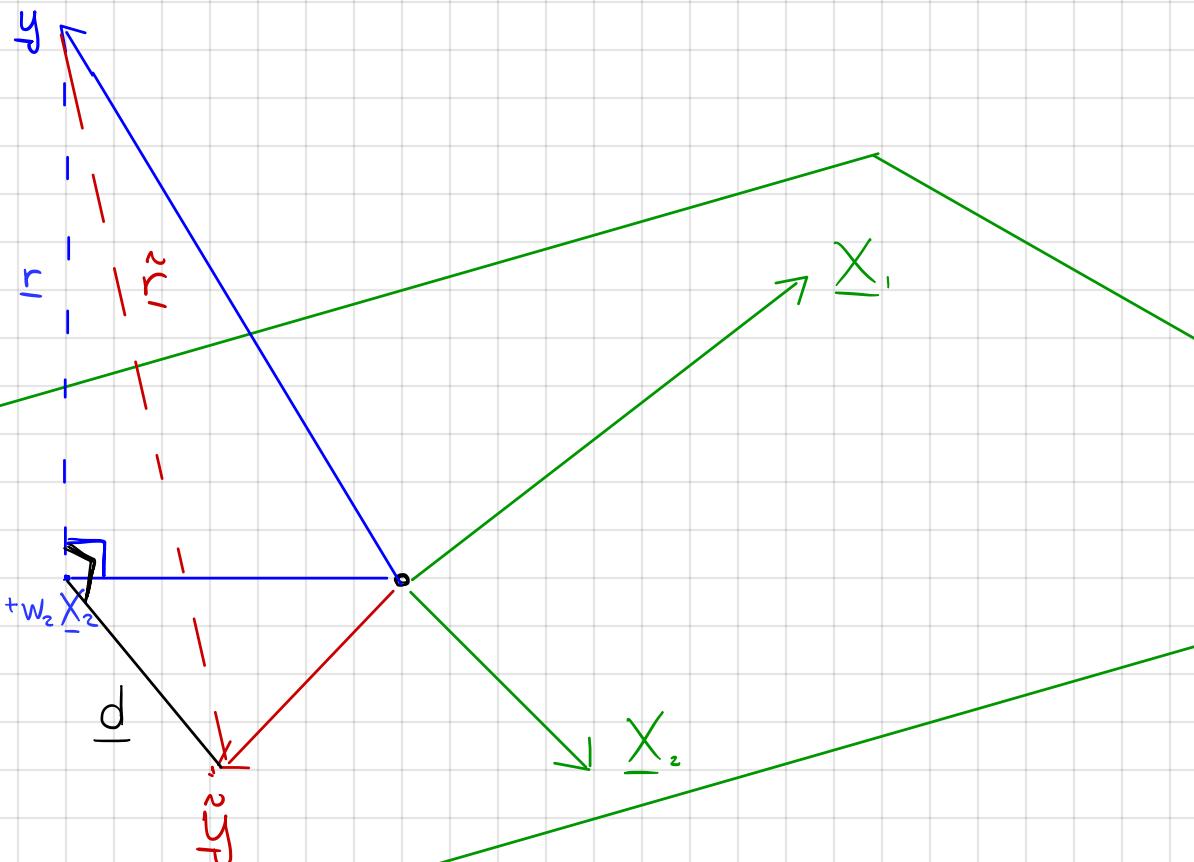
If $X = [\underline{X}_1, \underline{X}_2, \dots, \underline{X}_p]$, then $\text{range}(X) := \text{span}(\text{cols of } X) = \text{span}(\underline{X}_1, \dots, \underline{X}_p)$

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

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

Right angle because any other angle would correspond to a longer distance.
(bigger \tilde{r})



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

\underline{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 \underline{X} . y may not be in this space

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

$$= \underset{\underline{w}}{\operatorname{argmin}} \sum_{i=1}^n r_i^2(\underline{w})$$

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

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

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

$$\Rightarrow \underline{X}^\top \hat{\underline{r}} = \underline{0} \Leftarrow \text{vector of zeros}$$

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

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

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

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

- does it exist?
- is it unique?

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) $3x_1 + 2x_2 = 1$ $\Rightarrow x_2 = -3x_1$
 $3x_1 + x_2 = 0$ $-x_2 + 2x_2 = x_2 = 1 \Rightarrow x_1 = -\frac{1}{3}$ \Rightarrow one soln

b) $3x_1 + x_2 = 0$ ∞ solutions

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

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

e) $\begin{bmatrix} 3 & 2 \\ 3 & 1 \\ 2 & 2 \end{bmatrix} \begin{bmatrix} x_1 \\ x_2 \end{bmatrix} = \begin{bmatrix} 1 \\ 0 \\ 2 \end{bmatrix}$ $\Rightarrow 0$ solutions

f) $3x_1 + x_2 = 1$
 $6x_1 + 2x_2 = 2$ $\iff \begin{bmatrix} 3 & 1 \\ 6 & 2 \end{bmatrix} \begin{bmatrix} x_1 \\ x_2 \end{bmatrix} = \begin{bmatrix} 1 \\ 2 \end{bmatrix}$
 $\Rightarrow \infty$ solns rank-1 matrix!

Linear Independence

A collection of vectors $\underline{v}_1, \underline{v}_2, \dots, \underline{v}_p \in \mathbb{R}^n$ is linearly independent when
 $\sum_{i=1}^p \alpha_i \underline{v}_i = \underline{0}$ if and only if $\alpha_i = 0$ 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 $\alpha_1 = \alpha_2 = \alpha_3 = 0$

Ex $n=3$

$p=4$

$$\underline{v}_1 = \begin{bmatrix} 1 \\ 0 \\ 0 \end{bmatrix}, \quad \underline{v}_2 = \begin{bmatrix} 0 \\ 1 \\ 1 \end{bmatrix}, \quad \underline{v}_3 = \begin{bmatrix} 0 \\ 1 \\ 0 \end{bmatrix}, \quad \underline{v}_4 = \begin{bmatrix} 1 \\ 0 \\ 1 \end{bmatrix}$$

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

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

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

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

Recall the outer product representation of matrix product

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

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

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

Matrix Inverse

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

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

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

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

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

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

Recall from earlier:

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

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

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

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

$$\Rightarrow \underline{X}^\top \hat{\underline{r}} = \underline{0} \Leftarrow \text{vector of zeros}$$

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

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

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

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

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

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

$\underline{P}_{\underline{X}} := \underline{X}(\underline{X}^\top \underline{X})^{-1} \underline{X}^\top$ is called a Projection Matrix

because $\underline{P}_{\underline{X}} \underline{y}$ projects \underline{y} onto range(\underline{X})