

* Announcement — Gradescope submission.

Code submission. .PY

[HWO-A ← .pdf
HWO-B ← .pdf
HWO-code ← .py

	A	B
	✓	✓
		✓
	✓	✓

Lecture 3



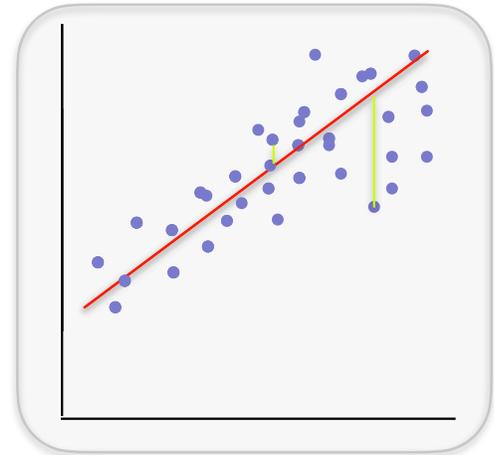
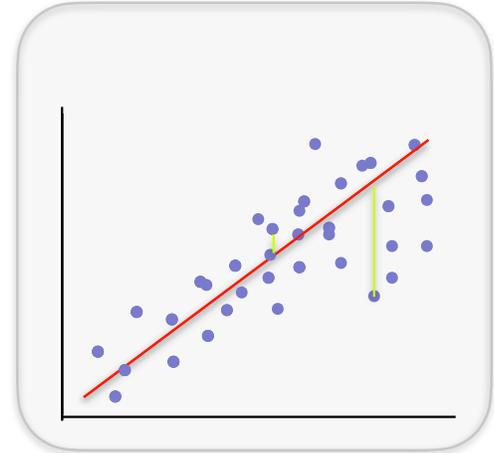
The regression problem in matrix notation

Linear model: $y_i = x_i^T w + \epsilon_i$

Least squares solution:

$$\begin{aligned} \mathbb{R}^d \ni \hat{w}_{LS} &= \arg \min_w \|\mathbf{y} - \mathbf{X}w\|_2^2 \\ &= (\mathbf{X}^T \mathbf{X})^{-1} \mathbf{X}^T \mathbf{y} \end{aligned}$$

What about an offset
(a.k.a intercept)?

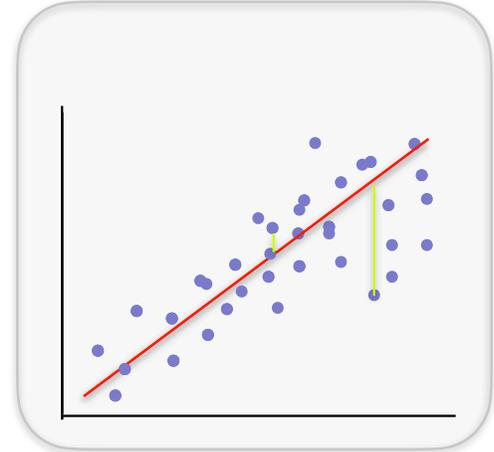


The regression problem in matrix notation

Linear model: $y_i = x_i^T w + \epsilon_i$

Least squares solution:

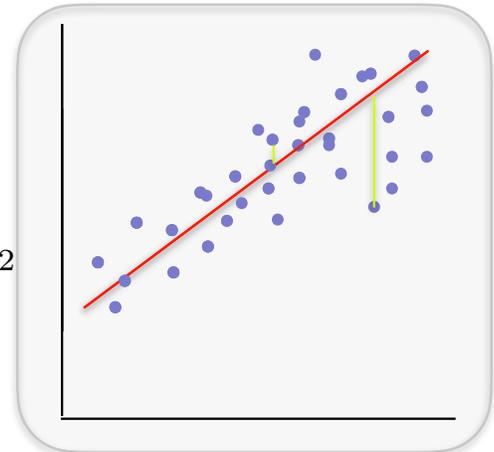
$$\begin{aligned}\hat{w}_{LS} &= \arg \min_w \|\mathbf{y} - \mathbf{X}w\|_2^2 \\ &= (\mathbf{X}^T \mathbf{X})^{-1} \mathbf{X}^T \mathbf{y}\end{aligned}$$



Affine model: $y_i = x_i^T w + \underset{\substack{\mathbb{F} \\ \mathbb{R}}}{b} + \epsilon_i$

Least squares solution:

$$\begin{aligned}\hat{w}_{LS}, \hat{b}_{LS} &= \arg \min_{w,b} \sum_{i=1}^n (y_i - (x_i^T w + b))^2 \\ &= \arg \min_{w,b} \|\mathbf{y} - (\mathbf{X}w + \mathbf{1}b)\|_2^2\end{aligned}$$



Dealing with an offset

$$\hat{w}_{LS}, \hat{b}_{LS} = \arg \min_{\substack{\mathbb{R}^d \ni w, b \\ \mathbb{R}}} \underbrace{\|y - (Xw + \mathbb{1}b)\|_2^2}_{\mathcal{L}(w, b) = (y - (Xw + \mathbb{1}b))^T (y - (Xw + \mathbb{1}b))} \leftarrow \text{Quad. fun. w.r.t. } w, b.$$

Vector $\begin{bmatrix} \vdots \\ \vdots \\ \vdots \end{bmatrix} \}^n$

Set gradient w.r.t. w and b to zero to find the minima:

$$\underbrace{\nabla_w}_{\mathbb{R}^d} \mathcal{L}(w, b) = 2 \cdot X^T (y - (Xw + \mathbb{1}b)) \Big|_{w, b=0}$$

$\underbrace{\begin{bmatrix} \vdots & \vdots & \vdots \end{bmatrix}}_n \left(\begin{bmatrix} \vdots \\ \vdots \\ \vdots \end{bmatrix} - \left(\begin{bmatrix} \vdots \\ \vdots \\ \vdots \end{bmatrix} + \begin{bmatrix} \mathbb{1} \\ \vdots \\ \vdots \end{bmatrix} b \right) \right)$

$$\underbrace{\nabla_b}_{\mathbb{R}} \mathcal{L}(w, b) = 2 \cdot \mathbb{1}^T (y - (Xw + \mathbb{1}b)) \Big|_{w, b=0}$$

Fact: $\nabla_w ((Aw + b)^T (Aw + b)) = 2A^T (Aw + b)$

$$(aw + b)^2 = 2a(aw + b) \quad \left| \begin{array}{l} \nabla_w (w^T w) = 2w, \\ \nabla_w (w^T A^T A w) = \end{array} \right. \rightarrow$$

$$X^T y = X^T X w + X^T \mathbb{1} b$$

$$\mathbb{1}^T y = \mathbb{1}^T X w + \mathbb{1}^T \mathbb{1} b$$

Linear in $\begin{matrix} \mathbb{R}^d \\ \mathbb{R}^1 \end{matrix}$ w, b

$d+1$ variables
 $d+1$ equations

Dealing with an offset

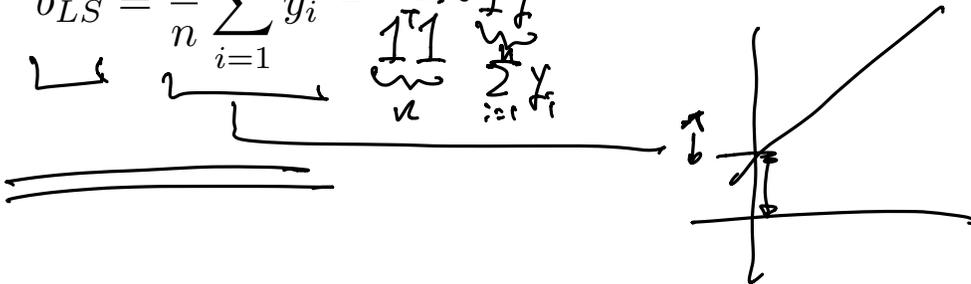
$$\hat{w}_{LS}, \hat{b}_{LS} = \arg \min_{w, b} \| \mathbf{y} - (\mathbf{X}w + \mathbf{1}b) \|_2^2$$

$$\begin{bmatrix} \sum_{i=1}^n x_{i,1} \\ \vdots \\ \sum_{i=1}^n x_{i,d} \end{bmatrix} = \begin{bmatrix} 0 \\ \vdots \\ 0 \end{bmatrix} \quad \begin{matrix} \mathbf{X}^T \mathbf{X} \hat{w}_{LS} + \hat{b}_{LS} \mathbf{X}^T \mathbf{1} = \mathbf{X}^T \mathbf{y} \\ \mathbf{1}^T \mathbf{X} \hat{w}_{LS} + \hat{b}_{LS} \mathbf{1}^T \mathbf{1} = \mathbf{1}^T \mathbf{y} \end{matrix} \rightarrow \begin{matrix} \hat{b} \cdot (\mathbf{1}^T \mathbf{1}) = \mathbf{1}^T \mathbf{y} \\ \hat{b} \cdot n = \sum_{i=1}^n y_i \\ \hat{b} = \frac{1}{n} \sum_{i=1}^n y_i \end{matrix}$$

If $\mathbf{X}^T \mathbf{1} = 0$, if the features have zero mean,

$$\hat{w}_{LS} = (\mathbf{X}^T \mathbf{X})^{-1} \mathbf{X}^T \mathbf{Y}$$

$$\hat{b}_{LS} = \frac{1}{n} \sum_{i=1}^n y_i = \frac{1}{n} \cdot \mathbf{1}^T \mathbf{y}$$



Dealing with an offset

$$\mu_j \triangleq \frac{1}{n} \sum_{i=1}^n X_{i,j}$$

μ_1 : average μ_2 : sum
latencies

$$\hat{w}_{LS}, \hat{b}_{LS} = \arg \min_{w,b} \underbrace{\| \mathbf{y} - (\mathbf{X}w + \mathbf{1}b) \|_2^2}_{\mathcal{L}(w,b)}$$

$$\begin{aligned} \mathbf{X}^T \mathbf{X} \hat{w}_{LS} + \hat{b}_{LS} \mathbf{X}^T \mathbf{1} &= \mathbf{X}^T \mathbf{y} & C &= \mu^T w + b \\ \mathbf{1}^T \mathbf{X} \hat{w}_{LS} + \hat{b}_{LS} \mathbf{1}^T \mathbf{1} &= \mathbf{1}^T \mathbf{y} & \mu &= \frac{1}{n} \mathbf{X}^T \mathbf{1} = \begin{bmatrix} \frac{1}{n} \sum_{i=1}^n X_{i,1} \\ \vdots \end{bmatrix} \end{aligned}$$

If $\mathbf{X}^T \mathbf{1} = 0$,

$$\hat{w}_{LS} = (\mathbf{X}^T \mathbf{X})^{-1} \mathbf{X}^T \mathbf{Y}$$

$$\hat{b}_{LS} = \frac{1}{n} \sum_{i=1}^n y_i$$

$\tilde{\mathbf{X}}, C$: new

In general, when $\mathbf{X}^T \mathbf{1} \neq 0$,

$$\begin{aligned} \mathcal{L}(w,b) &= \| \mathbf{y} - (\underbrace{\mathbf{X} - \mathbf{1}\mu^T}_{+}) w - \underbrace{\mathbf{1}\mu^T w - \mathbf{1}b}_{-} \|_2^2 \\ &= \| \mathbf{y} - \underbrace{\tilde{\mathbf{X}} \cdot w}_{\text{centered}} - \mathbf{1} \cdot C \|_2^2, \end{aligned}$$

$$\hat{w}_{LS} = (\tilde{\mathbf{X}}^T \tilde{\mathbf{X}})^{-1} \tilde{\mathbf{X}}^T \mathbf{Y}$$

$$\begin{aligned} \hat{C}_{LS} &= \frac{1}{n} \sum_{i=1}^n y_i \implies \hat{b}_{LS} = \frac{1}{n} \sum_{i=1}^n y_i - \mu^T \hat{w}_{LS} \\ &\parallel \\ \hat{b}_{LS} + \mu^T \hat{w}_{LS} \end{aligned}$$

Dealing with an offset

$$\hat{w}_{LS}, \hat{b}_{LS} = \arg \min_{w, b} \|\mathbf{y} - (\mathbf{X}w + \mathbf{1}b)\|_2^2$$

$$\mathbf{X}^T \mathbf{X} \hat{w}_{LS} + \hat{b}_{LS} \mathbf{X}^T \mathbf{1} = \mathbf{X}^T \mathbf{y}$$

$$\mathbf{1}^T \mathbf{X} \hat{w}_{LS} + \hat{b}_{LS} \mathbf{1}^T \mathbf{1} = \mathbf{1}^T \mathbf{y}$$

If $\mathbf{X}^T \mathbf{1} = 0$,

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

$$\hat{b}_{LS} = \frac{1}{n} \sum_{i=1}^n y_i$$

In general, when $\mathbf{X}^T \mathbf{1} \neq 0$,

$$\mu = \frac{1}{n} \mathbf{X}^T \mathbf{1}$$

$$\tilde{\mathbf{X}} = \mathbf{X} - \mathbf{1}\mu^T$$

$$\hat{w}_{LS} = (\tilde{\mathbf{X}}^T \tilde{\mathbf{X}})^{-1} \tilde{\mathbf{X}}^T \mathbf{y}$$

$$\hat{b}_{LS} = \frac{1}{n} \sum_{i=1}^n y_i - \mu^T \hat{w}_{LS}$$

Process

Decide on a **model**: $y_i = x_i^T w + b + \epsilon_i$

Choose a loss function - least squares

Pick the function which minimizes loss on data

$$\hat{w}_{LS}, \hat{b}_{LS} = \arg \min_{w, b} \sum_{i=1}^n (y_i - (x_i^T w + b))^2$$

Use function to make prediction on new examples

$$\hat{y}_{\text{new}} = x_{\text{new}}^T \hat{w}_{LS} + \hat{b}_{LS}$$

Why is least squares a good loss function?

$$\begin{aligned}\hat{w}_{LS} &= \arg \min_w \|\mathbf{y} - \mathbf{X}w\|_2^2 \\ &= (\mathbf{X}^T \mathbf{X})^{-1} \mathbf{X}^T \mathbf{y}\end{aligned}$$

Consider $y_i = x_i^T w + \epsilon_i$ where $\epsilon_i \stackrel{i.i.d.}{\sim} \mathcal{N}(0, \sigma^2)$

probabilistic
Generative
Model

$$\implies y_i \sim \mathcal{N}(x_i^T w, \sigma^2)$$

$$\implies P(y_i; x_i, w, \sigma) = \frac{1}{\sqrt{2\pi}\sigma} \cdot e^{-\frac{(x_i^T w - y_i)^2}{2\sigma^2}}$$

Why is least squares a good loss function?

Maximum Likelihood Estimator:

$$\begin{aligned}\hat{w}_{\text{MLE}} &= \arg \max_w \log P(\{y_i\}_{i=1}^n; \{x_i\}_{i=1}^n, w, \sigma) \\ &= \arg \max_w -n \log(\sigma\sqrt{2\pi}) + \sum_{i=1}^n -\frac{(y_i - x_i^T w)^2}{2\sigma^2} \\ &= \arg \max_w \sum_{i=1}^n -(y_i - x_i^T w)^2 \\ &= \arg \min_w \sum_{i=1}^n (y_i - x_i^T w)^2\end{aligned}$$

Why is least squares a good loss function?

Maximum Likelihood Estimator:

$$\begin{aligned}\hat{w}_{\text{MLE}} &= \arg \max_w \log P(\{y_i\}_{i=1}^n; \{x_i\}_{i=1}^n, w, \sigma) \\ &= \arg \max_w -n \log(\sigma \sqrt{2\pi}) + \sum_{i=1}^n -\frac{(y_i - x_i^T w)^2}{2\sigma^2} \\ &= \arg \min_w \sum_{i=1}^n (y_i - x_i^T w)^2\end{aligned}$$

Recall: $\hat{w}_{LS} = \arg \min_w \sum_{i=1}^n (y_i - x_i^T w)^2$

$$\hat{w}_{LS} = \hat{w}_{MLE} = (\mathbf{X}^T \mathbf{X})^{-1} \mathbf{X}^T \mathbf{Y}$$

Recap of linear regression

Data $\{(x_i, y_i)\}_{i=1}^n$

Minimize the loss (Empirical Risk Minimization)

Choose a loss

e.g., $(y_i - x_i^T w)^2$

Solve $\hat{w}_{\text{LS}} = \arg \min_w \sum_{i=1}^n (y_i - x_i^T w)^2$

Maximize the likelihood (MLE)

Choose a Hypothesis class

e.g., $y_i = x_i^T w + \epsilon_i, \quad \epsilon_i \sim \mathcal{N}(0, \sigma^2)$

Maximize the likelihood,

$$\hat{w}_{\text{MLE}} = \arg \max_w \left\{ -n \log(\sigma \sqrt{2\pi}) - \sum_{i=1}^n \frac{(y_i - x_i^T w)^2}{2\sigma^2} \right\}$$

Analysis of **Error** under additive Gaussian noise

$$\text{if } y_i = x_i^T w + \epsilon_i \quad \text{and} \quad \epsilon_i \stackrel{i.i.d.}{\sim} \mathcal{N}(0, \sigma^2) \quad \mathbf{Y} = \mathbf{X}w + \epsilon$$

$$\begin{aligned} \hat{w}_{MLE} &= (\mathbf{X}^T \mathbf{X})^{-1} \mathbf{X}^T \mathbf{Y} \\ &= (\mathbf{X}^T \mathbf{X})^{-1} \mathbf{X}^T (\mathbf{X}w + \epsilon) \\ &= w + (\mathbf{X}^T \mathbf{X})^{-1} \mathbf{X}^T \epsilon \end{aligned}$$

Handwritten annotations:
- A bracket labeled n is drawn around the ϵ term in the second equation.
- An arrow labeled $\mathcal{N}(0, \sigma^2)$ points to the ϵ term in the second equation.
- A circle around the ϵ term in the third equation has an arrow pointing to the word "Random".

Maximum Likelihood Estimator is unbiased:

$$\begin{aligned} \text{Bias}(\hat{w}_{MLE}) &= \mathbb{E}[\hat{w}_{MLE}] - w \\ &= \mathbb{E}[w + (\mathbf{X}^T \mathbf{X})^{-1} \mathbf{X}^T \epsilon] - w \\ &= w + (\mathbf{X}^T \mathbf{X})^{-1} \mathbf{X}^T \mathbb{E}[\epsilon] - w \\ &= w + (\mathbf{X}^T \mathbf{X})^{-1} \mathbf{X}^T \mathbf{0} - w \\ &= w - w = 0 \end{aligned}$$

Handwritten annotations:
- A vertical line with a double arrow is drawn under the $\mathbb{E}[\epsilon]$ term, with a "0" written below it.
- A vertical line with a double arrow is drawn under the $\mathbf{0}$ term in the final step, with a "0" written below it.

Analysis of **Error** under additive Gaussian noise

if $y_i = x_i^T w + \epsilon_i$ and $\epsilon_i \stackrel{i.i.d.}{\sim} \mathcal{N}(0, \sigma^2)$ $\mathbf{Y} = \mathbf{X}w + \epsilon$

$$\begin{aligned} \hat{w}_{MLE} &= (\mathbf{X}^T \mathbf{X})^{-1} \mathbf{X}^T \mathbf{Y} \\ &= (\mathbf{X}^T \mathbf{X})^{-1} \mathbf{X}^T (\mathbf{X}w + \epsilon) \\ &= w + (\mathbf{X}^T \mathbf{X})^{-1} \mathbf{X}^T \epsilon \end{aligned} \quad \Bigg]$$

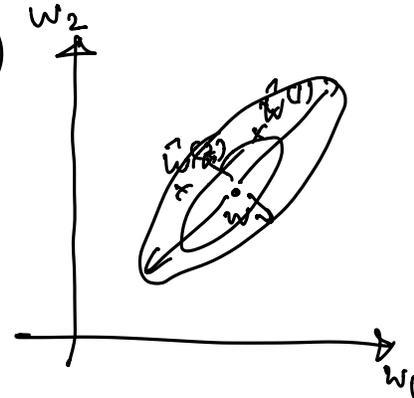
Covariance is: $\mathbb{E}[(\hat{w} - w)(\hat{w} - w)^T]$

$$\begin{aligned} &= \mathbb{E}\left[(\mathbf{X}^T \mathbf{X})^{-1} \mathbf{X}^T \epsilon \cdot \epsilon^T \mathbf{X} (\mathbf{X}^T \mathbf{X})^{-1}\right] \\ &= (\mathbf{X}^T \mathbf{X})^{-1} \mathbf{X}^T \underbrace{\mathbb{E}[\epsilon \epsilon^T]}_{\sigma^2 \mathbf{I}} \mathbf{X} (\mathbf{X}^T \mathbf{X})^{-1} = \sigma^2 \underbrace{(\mathbf{X}^T \mathbf{X})^{-1} \mathbf{X}^T \mathbf{X} (\mathbf{X}^T \mathbf{X})^{-1}}_{\square \cdot \square^T = \mathbf{I}} \\ &= \sigma^2 (\mathbf{X}^T \mathbf{X})^{-1} \end{aligned}$$

Analysis of **Error** under additive Gaussian noise

$$\text{if } y_i = x_i^T w + \epsilon_i \quad \text{and} \quad \epsilon_i \stackrel{i.i.d.}{\sim} \mathcal{N}(0, \sigma^2) \quad \mathbf{Y} = \mathbf{X}w + \epsilon$$

$$\begin{aligned} \hat{w}_{MLE} &= (\mathbf{X}^T \mathbf{X})^{-1} \mathbf{X}^T \mathbf{Y} \\ &= (\mathbf{X}^T \mathbf{X})^{-1} \mathbf{X}^T (\mathbf{X}w + \epsilon) \\ &= w + \underbrace{(\mathbf{X}^T \mathbf{X})^{-1} \mathbf{X}^T}_{\text{matrix}} \epsilon \end{aligned}$$



$$\mathbb{E}[\hat{w}_{MLE}] = w$$

$$\text{Cov}(\hat{w}_{MLE}) = \mathbb{E}[(\hat{w} - \mathbb{E}[\hat{w}])(\hat{w} - \mathbb{E}[\hat{w}])^T] = (\mathbf{X}^T \mathbf{X})^{-1} \cdot \sigma^2$$

$$\hat{w}_{MLE} \sim \mathcal{N}(w, \underbrace{(\mathbf{X}^T \mathbf{X})^{-1} \cdot \sigma^2}_{\text{covariance matrix}})$$

Questions?

① What if $X^T X$ has rank $< d$.
 $d \times d$

$n > d$ & X_i 's in general positions.

$$v^T (X^T X) \cdot v = 0.$$

\hat{w}_{MLE}

~~$(X^T X)^+ \cdot X^T \cdot y + \text{cov}$~~



all in this line solve with loss.