

Lecture 3: Linear regression (continued)



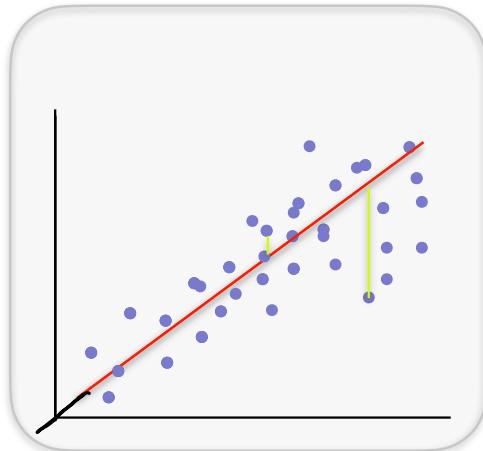
W

The regression problem in matrix notation

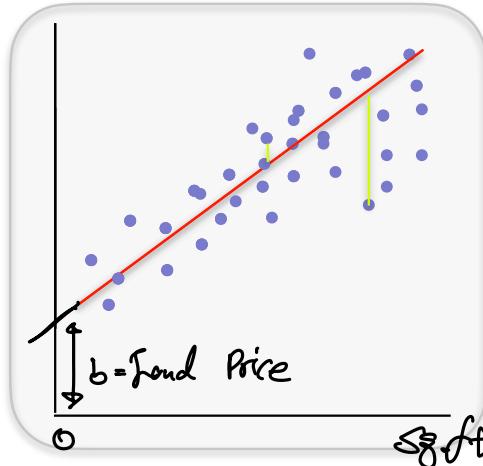
Linear model: $y_i = x_i^T w + \epsilon_i$
 $x_i, w \in \mathbb{R}^d$

Least squares solution:

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



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

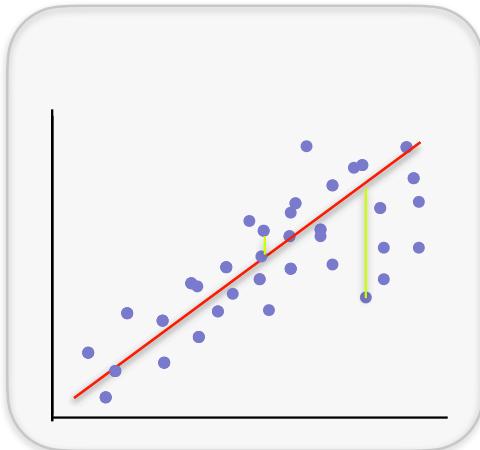


The regression problem in matrix notation

Linear model: $y_i = x_i^T w + \epsilon_i$
 $x_i, w \in \mathbb{R}^d$

Least squares solution:

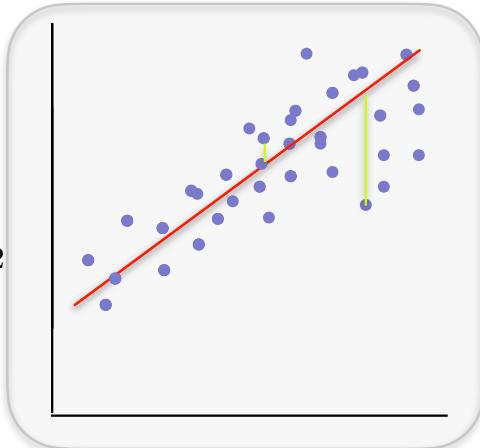
$$\begin{aligned}\hat{w}_{LS} &= \arg \min_w \|\mathbf{y} - \mathbf{X}w\|_2^2 \\ &= (\mathbf{X}^T \mathbf{X})^{-1} \mathbf{X}^T \mathbf{y}\end{aligned}$$



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

Least squares solution:

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



Dealing with an offset

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

=

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

$$\nabla_w \mathcal{L}(w,b) = \mathbf{X}^\top(\mathbf{y} - \mathbf{X}w - \mathbf{1}b) \quad \text{ER}^d \quad \left| \begin{array}{l} \nabla_b \mathcal{L}(w,b) = \mathbf{1}^\top(\mathbf{y} - \mathbf{X}w - \mathbf{1}b) \quad \text{ER} \\ \rightarrow \mathbf{X}^\top \mathbf{y} = \mathbf{X}^\top \mathbf{X} \hat{w} + \mathbf{1}^\top \mathbf{1} \cdot \hat{b} \quad \text{GR} \\ \left[\begin{array}{c} \mathbf{X}^\top \\ \mathbf{1}^\top \end{array} \right] \cdot \mathbf{y} = \left[\begin{array}{c} \mathbf{X}^\top \\ \mathbf{1}^\top \end{array} \right] \left[\begin{array}{cc} \mathbf{X} & \mathbf{1} \end{array} \right] \left[\begin{array}{c} \hat{w} \\ \hat{b} \end{array} \right] \end{array} \right.$$

* Series of linear equations with
 $d+1$ -dim variables (\hat{w}, \hat{b})
 $d+1$ equations.

• How do you solve this?

$$\left[\begin{array}{c} \mathbf{X}^\top \\ \mathbf{1}^\top \end{array} \right] \cdot \mathbf{y} = \left[\begin{array}{c} \mathbf{X}^\top \\ \mathbf{1}^\top \end{array} \right] \left[\begin{array}{c} \mathbf{X} \\ \mathbf{1} \end{array} \right] \left[\begin{array}{c} \hat{w} \\ \hat{b} \end{array} \right]$$

Schur complement $\left(\left[\begin{array}{c} \mathbf{X}^\top \\ \mathbf{1}^\top \end{array} \right] \left[\begin{array}{c} \mathbf{X} \\ \mathbf{1} \end{array} \right] \right)^{-1} \left[\begin{array}{c} \mathbf{X}^\top \\ \mathbf{1}^\top \end{array} \right] \cdot \mathbf{y} = \left[\begin{array}{c} \hat{w} \\ \hat{b} \end{array} \right]$

Dealing with an offset

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

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

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

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

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

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

i-th sample
 $\frac{1}{n} \mathbf{X}^T \mathbf{1} = \frac{1}{n} \begin{bmatrix} x_{11} & \dots & x_{1j} & \dots & x_{1n} \end{bmatrix} \mathbf{1}^T$
j-th feature
 $= \left[\frac{1}{n} \sum_{i=1}^n x_{i,j} \right] = \frac{1}{n} \sum_{i=1}^n x_i = \mu$
2nd feature
 $= \mu$
 $\mu_1: \text{Average Sq. ff.}$
 $\mu_2: \text{Average bathroom ff.}$
1st feature

Dealing with an offset

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

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

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

$$c = b + \mu^T w$$

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

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

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

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

\hat{X}, c

In general, when $\mathbf{X}^T \mathbf{1} \neq 0$,
 $\mathcal{L}(w,b) = \|\mathbf{y} - \mathbf{X}w + \mathbf{1}\mu^T w - \mathbf{1}\mu^T w - \mathbf{1}b\|_2^2$

$$\mathcal{L}(w,c) = \|\mathbf{y} - (\mathbf{X} - \mathbf{1}\mu^T)w - \mathbf{1} \cdot c\|_2^2$$

$\xrightarrow{\text{Zero-Mean } \hat{X}}$

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

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

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

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

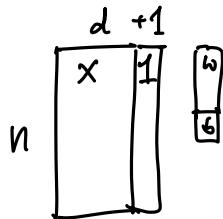
$$\hat{y}_{\text{new}} = x_{\text{new}}^T \hat{w}_{LS} + \hat{b}_{LS}$$

Another way of dealing with an offset

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

reparametrize the problem as $\bar{\mathbf{X}} = [\mathbf{X}, \mathbf{1}]$ and $\bar{w} = \begin{bmatrix} w \\ b \end{bmatrix}$

$$\bar{\mathbf{X}} \bar{w} = \begin{bmatrix} \mathbf{X} & \mathbf{1} \end{bmatrix} \cdot \begin{bmatrix} w \\ b \end{bmatrix} = \mathbf{X} \cdot w + \mathbf{1} \cdot b$$



$$\hat{\bar{w}}_{LS} = \arg \min \|\mathbf{y} - \bar{\mathbf{X}} \cdot \bar{w}\|_2^2 = (\bar{\mathbf{X}}^\top \bar{\mathbf{X}})^{-1} \cdot \bar{\mathbf{X}}^\top \cdot \mathbf{y}$$
$$\begin{bmatrix} \hat{w} \\ \hat{b} \end{bmatrix} = \left(\begin{bmatrix} \mathbf{X}^\top \\ \mathbf{1}^\top \end{bmatrix} \begin{bmatrix} \mathbf{X} & \mathbf{1} \end{bmatrix} \right)^{-1} \cdot \begin{bmatrix} \mathbf{X} & \mathbf{1} \end{bmatrix} \cdot \mathbf{y}$$

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 $\epsilon_i \stackrel{i.i.d.}{\sim} \mathcal{N}(0, \sigma^2)$

$$\begin{aligned}\Rightarrow y_i &\sim \mathcal{N}(x_i^T w, \sigma^2) \\ \Rightarrow P(y_i; x_i, w, \sigma) &= \frac{1}{\sqrt{2\pi\sigma^2}} e^{-\frac{(y_i - x_i^T w)^2}{2\sigma^2}}\end{aligned}$$

Why is least squares a good loss function?

Maximum Likelihood Estimator:

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

Why is least squares a good loss function?

Maximum Likelihood Estimator:

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

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

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

* Gaussian noise \iff ℓ_2 -loss minimization.

Recap of linear regression

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

**Minimize the loss
(Empirical Risk Minimization)**

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

there are other choices

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

there are other hypotheses

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

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

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

↑ Random

Maximum Likelihood Estimator is unbiased:

$$\begin{aligned}\mathbb{E}[\hat{w}_{MLE}] &= w^* + \mathbb{E}[(\mathbf{X}^T \mathbf{X})^{-1} \mathbf{X}^T \epsilon] \quad \leftarrow \text{Linearity of expectation} \\ &\stackrel{?}{=} w^* + (\mathbf{X}^T \mathbf{X})^{-1} \mathbf{X}^T \cdot \mathbb{E}[\epsilon] \quad \leftarrow \mathbb{E}[a \cdot \epsilon + b] = a \mathbb{E}[\epsilon] + b \\ &= w^*\end{aligned}$$

Analysis of Error under additive Gaussian noise

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

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

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

linearity of expectation

$$\begin{aligned}&= \mathbb{E}[(x^T x)^{-1} x^T \epsilon \cdot \epsilon^T x (x^T x)^{-1}] \\ &= (x^T x)^{-1} x^T \mathbb{E}[\epsilon \epsilon^T] x (x^T x)^{-1} \\ &= (x^T x)^{-1} x^T \sigma^2 I \cdot x (x^T x)^{-1} \\ &= \sigma^2 (x^T x)^{-1} x^T x (x^T x)^{-1} \\ &= \sigma^2 (x^T x)^{-1}\end{aligned}$$

$$\begin{aligned}\mathbb{E}[\epsilon \epsilon^T] &= M \\ M_{ij} &= \mathbb{E}[\epsilon_i \epsilon_j] \\ &= \begin{cases} \sigma^2 & \text{if } i=j \\ 0 & \text{else} \end{cases} \\ M &= \begin{bmatrix} \sigma^2 & & & \\ & \ddots & & \\ & & \ddots & \\ & & & 0 \end{bmatrix} = \sigma^2 I.\end{aligned}$$

Analysis of Error under additive Gaussian noise

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

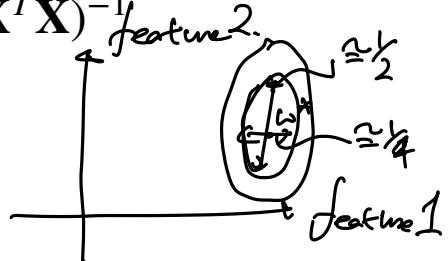
$$\begin{aligned}\hat{w}_{MLE} &= (\mathbf{X}^T \mathbf{X})^{-1} \mathbf{X}^T \mathbf{Y} \\ &= (\mathbf{X}^T \mathbf{X})^{-1} \mathbf{X}^T (\mathbf{X}w + \epsilon) \\ &= w + (\mathbf{X}^T \mathbf{X})^{-1} \mathbf{X}^T \epsilon\end{aligned}$$

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

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

$$\hat{w}_{MLE} \sim \mathcal{N}(w, \sigma^2 (\mathbf{X}^T \mathbf{X})^{-1}) = \mathcal{N}(w, \sigma^2 \begin{bmatrix} \frac{1}{4} & 0 \\ 0 & \frac{1}{2} \end{bmatrix})$$

$$\text{If } \mathbf{X} = \begin{bmatrix} 1 & 0 \\ 0 & 1 \\ -1 & 0 \\ 0 & -1 \end{bmatrix} \text{ then } \mathbf{X}^T \mathbf{X} = \begin{bmatrix} 4 & 0 \\ 0 & 2 \end{bmatrix} \quad (\mathbf{X}^T \mathbf{X})^{-1} = \begin{bmatrix} \frac{1}{4} & 0 \\ 0 & \frac{1}{2} \end{bmatrix}$$



Questions?
