

STAT 391

2/25/2025

# Lecture 15

Linear Regression  
- Double Descent

Q2 Thu 2/27

# Lecture Notes Regression

Marina Meilă  
[mmp@stat.washington.edu](mailto:mmp@stat.washington.edu)

Department of Statistics  
University of Washington



Prediction problems ✓

Linear regression

← Examples  
 $r^2$

(Linear regression for non-linear  $f$ )

Is model  $\approx$  true?

Reading: Ch.

# Linear regression

## Model

- ▶  $y^i = \beta_0 + \beta_1 x_1^i + \epsilon^i$  for  $x^{1:n} \in \mathbb{R}$  (univariate regression)
  - ▶  $y^i = \beta_0 + \beta_1 x_1^i + \beta_2 x_2^i + \dots + \beta_d x_d^i + \epsilon^i$  for  $x^{1:n} \in \mathbb{R}^d$  (multivariate regression)
- $y \in \mathbb{R}$  is output/response  
 $x_j^{1:n}$  for  $j = 1 : d$  are input(s)/covariates/features/attributes/...  
 $\beta_{1:d}$  are regression coefficients,  $\beta_0$  is intercept  
 $\epsilon^{1:n} \in \mathbb{R}$  is noise,  $\epsilon^{1:n} \sim N(0, \sigma^2)$  i.i.d.

Data Model

$$(x^1, y^1), \dots, (x^n, y^n)$$

$$y = \beta_0 + \beta^T x + \epsilon \text{ noise } \sim N(0, \sigma^2)$$

$$= [\beta_0 \ \boxed{\beta_1 \ \dots \ \beta_d}] \begin{bmatrix} 1 \\ x_1 \\ \vdots \\ x_d \end{bmatrix} + \epsilon$$

Wanted

params:  $\underset{=}{{\beta_0:q}}$   $\sigma^2$  <sup>ML</sup>

$x$  = covariates  
regressors  
attributes  
features  
inputs

$y$  = output  
dependent  
variable

## Linear regression model in matrix form

For a single data point  $(x^i, y^i)$

$$\beta = \begin{bmatrix} \beta_1 \\ \beta_2 \\ \dots \\ \beta_d \end{bmatrix} \in \mathbb{R}^d \quad x^i = \begin{bmatrix} x_1^i \\ x_2^i \\ \dots \\ x_d^i \end{bmatrix} \in \mathbb{R}^d \quad (1)$$

Then,

$$y^i = \beta_0 + (x^i)^T \beta + \epsilon^i \quad (2)$$

For all  $\mathcal{D}$

$$y = \begin{bmatrix} y^1 \\ y^2 \\ \dots \\ y^n \end{bmatrix} \in \mathbb{R}^n \quad X = \begin{bmatrix} (x^1)^T \\ (x^2)^T \\ \dots \\ (x^n)^T \end{bmatrix} \in \mathbb{R}^{n \times d} \quad \epsilon = \begin{bmatrix} \epsilon^1 \\ \epsilon^2 \\ \dots \\ \epsilon^n \end{bmatrix} \quad (3)$$

$$y = \beta_0 \mathbf{1} + X\beta + \epsilon \quad \text{with } Cov(\epsilon) = \sigma^2 I_d \quad (4)$$

## Estimation of $\sigma^2$

- ▶ Let  $\hat{\epsilon}^i = y^i - (x^i)^T \beta^{ML}$ , for  $i = 1 : n$
- ▶  $\hat{\epsilon}^i$  called the **residuals**
- ▶ If we plug in  $\beta^{ML}$  in the log-likelihood equation (12) we obtain

$$I(\sigma^2, \beta^{ML}) = -\frac{n}{2} \ln \sigma^2 - \frac{1}{2\sigma^2} \sum_{i=1}^n (\hat{\epsilon}^i)^2 + \text{constant} \quad (21)$$

- ▶ Maximizing this expression w.r.t.  $\sigma^2$  gives

$$(\sigma^2)^{ML} = \frac{1}{n} \sum_{i=1}^n (\hat{\epsilon}^i)^2 \quad (22)$$

- ▶ The predictor is  $f(x) = x^T \beta^{ML}$
- ▶ Hence, for a new  $x \in \mathbb{R}^d$ , our guess of  $y$  is  $\hat{y} = f(x)$

## Maximizing the log-likelihood (w.r.t $\beta$ ) w.r.t $\sigma^2$

- For simplicity, let  $\beta_0 = 0$ ; hence  $y^i = \beta^T x^i + \epsilon^i$
- log-likelihood

$$I(\beta_{1:d}, \sigma^2) = -\frac{n}{2} \ln \sigma^2 - \frac{1}{2\sigma^2} \sum_{i=1}^n (y^i - \beta_0 - \beta^T x^i)^2 + \text{constant} \quad (12)$$

$\rightarrow$  indep of  $\sigma^2$

- For any  $\sigma^2$ ,

$$\underset{\beta}{\operatorname{argmax}} I(\sigma^2, \beta) = \underset{\beta}{\operatorname{argmin}} \sum_{i=1}^n (y^i - \beta_0 - \beta^T x^i)^2 \quad (13)$$

### a Least Squares Problem

- In matrix form  $\min_{\beta} \|y - X\beta\|^2$

- Solution

$$\beta^{ML} = (X^T X)^{-1} X^T y \quad (14)$$

with  $(X^T X)^{-1} X^T \equiv X^\dagger$  the **pseudoinverse** of  $X$

$\beta^{ML}$  is **linear** in  $y!$

$\beta^{ML}$  is **non-linear**,  $X$

ML for  $\sigma^2$   $\max_{\sigma^2} -\frac{n}{2} \ln \sigma^2 - \frac{1}{2\sigma^2} [C]$

$$\max_{\sigma^2} -\frac{n}{2} \ln \frac{\sigma^2}{n} - \frac{1}{2\sigma^2} [C] \quad \text{Calculus} \checkmark$$

$$\frac{\partial L}{\partial (\sigma^2)} = -\frac{n}{2} \cdot \frac{1}{\sigma^2} + \frac{1}{2} \cdot \frac{1}{(\sigma^2)^2} \stackrel{C}{=} 0 \Rightarrow \frac{C}{\sigma^2} = n \Rightarrow (\sigma^2)^M = \frac{1}{n} \sum_{i=1}^n (y^i - (x^i) f^{\text{ML}})^2$$

$\hat{\varepsilon}^i$  = estimate of noise  $\varepsilon^i$

Best predictor:  $\underset{\text{(linear)}}{\operatorname{argmin}} \sum \text{err}^2$

avg. of  
(errors)<sup>2</sup>

"LEAST SQUARES"  
least sum of  
Squared Errors

$$= \frac{1}{n} \sum_{i=1}^n (\hat{\varepsilon}^i)^2 \approx \text{Variance !!}$$

$$\text{Ex: } s = \sum_{i=1}^n \hat{\varepsilon}^i \cdot \frac{1}{n} = 0$$

Model  $\varepsilon^i \sim N(0, \sigma^2)$

ML Estimation (Training)  $\uparrow$   
Prediction  $\downarrow$   $f^{\text{ML}}, (\sigma^2)^{\text{ML}}$   
 (Testing)

Predictor  $f(x) = (\beta^{\text{ML}})^T x$  linear predictor  
 new  $x$   $\rightarrow \hat{y}$  prediction of  $y$  given  $x$

## Statistical properties of $\beta^{ML}$

- ▶ Assume the true model is  $y^i = \beta^T x^i + \epsilon^i$  with  $\epsilon \sim N(0, \sigma^2)$ .
- ▶ Here  $\beta, \sigma^2$  are **true parameters** and we assume we know them.
- ▶  $\beta^{ML}$  is a random variable. Let us calculate its mean and standard deviation.
- ▶ **Expectation** of  $\beta^{ML}$

*estimated*

$$E[\beta^{ML}] = E[X^\dagger y] = E[X^\dagger(X\beta + \epsilon)] \quad (15)$$

$$= E[X^\dagger X]\beta + E[X^\dagger \epsilon] \quad (16)$$

$$= \underbrace{X^\dagger X}_{I_d} \beta + \underbrace{X^\dagger E[\epsilon]}_0 = \beta \quad (17)$$

Assume Model True

Hence,  $E[\beta^{ML}] = \beta$  and we say that  $\beta^{ML}$  is **unbiased**

$$\beta^{ML} \sim N(\beta, \Sigma_\beta)$$

- ▶ **Covariance** of  $\beta^{ML}$
- ▶ By a similar (but longer) calculation, we obtain that

$$\text{Cov}(\beta^{ML}) = \sigma^2 (X X^\dagger)^{-1} \in \mathbb{R}^{d \times d} \quad (18)$$

- ▶ In fact,  $\beta^{ML}$  has a Normal distribution. Remember from (15)

$$\beta^{ML} = \beta + X^\dagger \epsilon. \quad (19)$$

This is a linear transformation of  $\epsilon$ , a Gaussian variable. Therefore,

$$\beta^{ML} \sim N(\beta, \sigma^2 (X X^\dagger)^{-1}). \quad (20)$$

- ▶ This is a **multivariate Normal** distribution over  $\mathbb{R}^d$ ,

Cov  $\hat{\beta}^{\text{ML}}$

$$\Sigma_{\beta} = \begin{bmatrix} \text{Var } \hat{\beta}_0^{\text{ML}} & \text{Cov}(\hat{\beta}_i^{\text{ML}}, \hat{\beta}_j^{\text{ML}}) \\ \text{Cov}(\hat{\beta}_i^{\text{ML}}, \hat{\beta}_j^{\text{ML}}) & \text{Var } \hat{\beta}_d^{\text{ML}} \end{bmatrix} \quad i, j = 0 : d$$

$$\Sigma_{\beta} = \sigma^2 (X^T X)^{-1} = \frac{\sigma^2}{n} \cdot \frac{1}{\hat{\sigma}^2_X} \quad \begin{array}{l} \text{noise} \\ \text{Input data} \end{array}$$

$d=1, \beta_0=0 \Rightarrow y = \beta x + \varepsilon \Rightarrow x, \beta \in (-\infty, \infty)$

Assume  $\sum x^i = 0 \Leftrightarrow \frac{1}{n} \sum x^i = 0$

$$\sigma_{\hat{\beta}}^2 = \frac{1}{n} \frac{\sigma^2}{\hat{\sigma}^2_X} \propto \frac{\text{noise}}{X \text{ data} \cdot \text{"spread"}}$$

$\rightarrow 0$  with  $n$

$$\begin{bmatrix} x^1 \\ \vdots \\ x^n \end{bmatrix} = X \Rightarrow X^T X = \sum_{i=1}^n (x^i)^2 = n \cdot \text{Var } x^i \quad \begin{array}{l} \text{sample} \\ \text{variance} \end{array}$$

$$\frac{\hat{\sigma}^2_X}{\sigma^2_{\text{noise}}} = \text{SNR} = \frac{\text{Signal to Noise Ratio}}{\text{Noise Ratio}}$$

$\text{SNR} \uparrow \Rightarrow \sigma_{\hat{\beta}}^2 \downarrow$

When model  $y = \beta_0 + \beta_1 x_1 + \dots + \beta_d x_d + \varepsilon$  NOT true

①  $\hat{\beta}^{\text{ML}}$  Estimation

Lin reg how used?

② Test it →

Interpret coefficients  $\beta_0, \beta_1, \dots$   
(Hypothesis testing  
Confidence Intervals)

"Interpolation"

inference

$\beta_j \approx 0 \implies x_j$  not important for  $y$

$\beta_j > 0$

$x_j$  has significant effect on  $y$

$\beta_j < 0$

## diagnosis

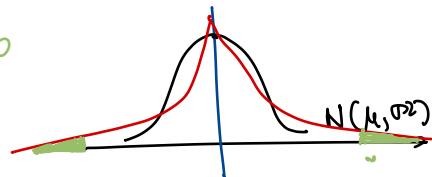
1. Model  $\approx$  true?
2. approx true but outlier

+ wrong date  
not from P data source  
from P when P has heavy tails

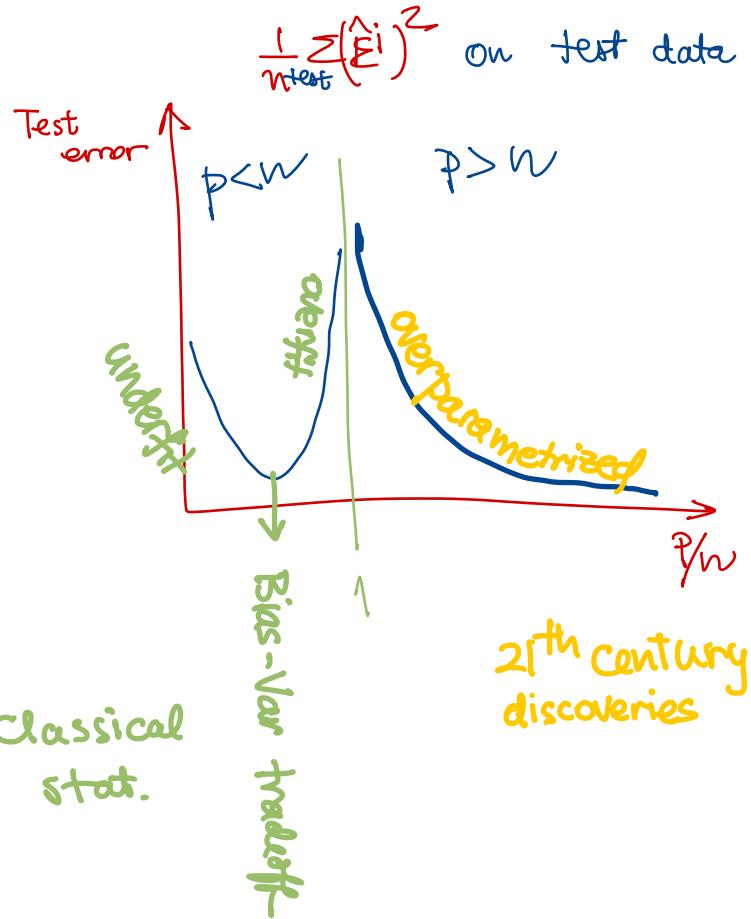
HAVE LARGE  
~~EFFECT~~  
ON BML !!  
INFLUENCE

statistically  
 $P(\text{out} | X) \approx 0$   
by your model

$P(\text{tail}) \neq 0$



# Bias-Variance Revisited



$$n = 300 \quad n_{\text{test}} = 500$$

$$\begin{array}{l} p < n \\ p > n \\ y = \alpha^T x + \varepsilon \end{array}$$

True

$$x \in \mathbb{R}^{100}$$

Neural net

$$x \in \mathbb{R}^P$$

$$\frac{p}{n} \in (0, 5)$$

Model

$$y = \beta^T x + \varepsilon$$