

Lecture 14:

supervised learning

Regression $\rightarrow y \in \mathbb{R}$

Classification

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

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

$$y_i \in \mathbb{R}$$

$$\rightarrow y_i \approx f(x_i)$$

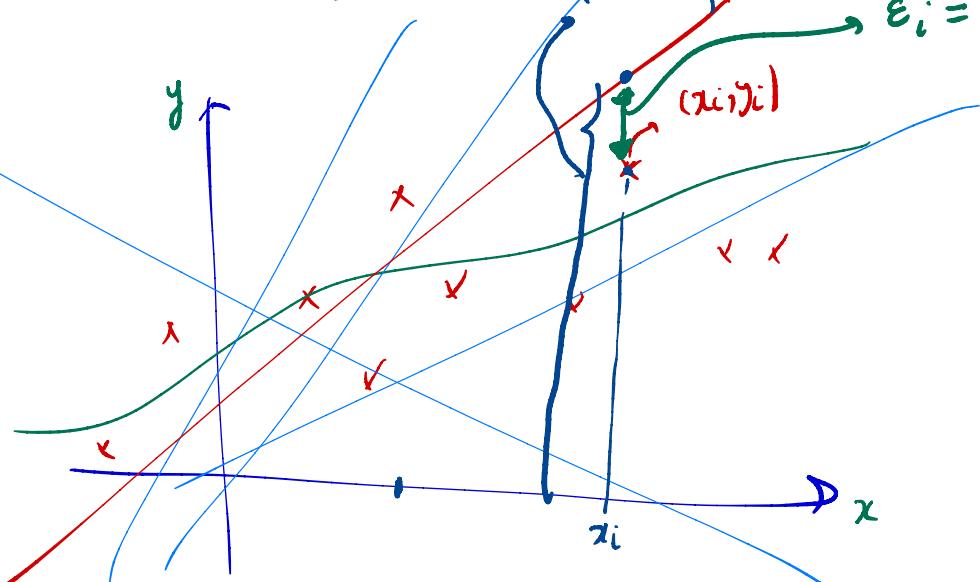
will be used for prediction on newly
unseen data points

$$P=1$$

$$\begin{cases} x_i \in \mathbb{R} \\ y_i \in \mathbb{R} \end{cases}$$

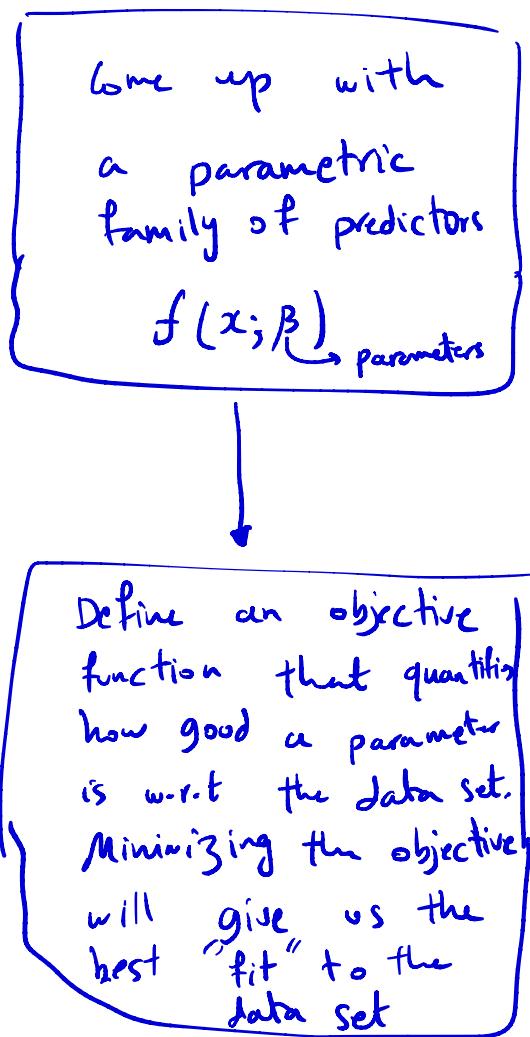
$$f(x; \beta)$$

$$\epsilon_i = y_i - f(x_i; \beta)$$



The simplest parametric family: $y = f(x; \beta) = \hat{\beta}_0 + \hat{\beta}_1 x$

our general
methodology:

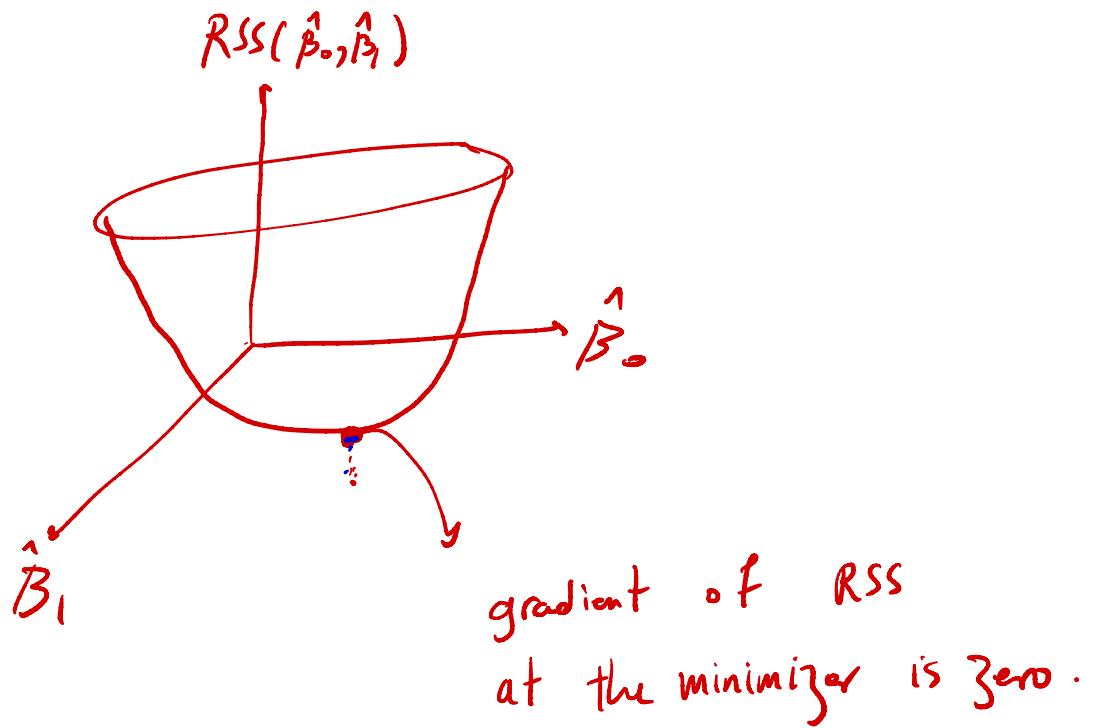


objective for the linear model:

$$\min_{\hat{\beta}_0, \hat{\beta}_1} \sum_{i=1}^n \epsilon_i^2 \stackrel{\Delta}{=} \min_{\hat{\beta}_0, \hat{\beta}_1} \sum_{i=1}^n (y_i - (\hat{\beta}_0 + \hat{\beta}_1 x_i))^2$$

$\underbrace{\qquad\qquad\qquad}_{\text{RSS } (\hat{\beta}_0, \hat{\beta}_1)}$

Residual Sum of Squares



These two equations
Should hold at
the minimizer.

$$\left. \begin{aligned} \frac{d \text{RSS}(\hat{\beta}_0, \hat{\beta}_1)}{d \hat{\beta}_0} &= 0 \\ \frac{d \text{RSS}(\hat{\beta}_0, \hat{\beta}_1)}{d \hat{\beta}_1} &= 0 \end{aligned} \right\}$$

$$\text{RSS}(\hat{\beta}_0, \hat{\beta}_1) = \sum_{i=1}^n \left(y_i - \hat{\beta}_0 + \hat{\beta}_1 x_i \right)^2$$

$$\frac{d}{d\hat{\beta}_0} RSS(\hat{\beta}_0, \hat{\beta}_1) \quad (1)$$

$$= -2 \sum_{i=1}^n (y_i - \hat{\beta}_0 - \hat{\beta}_1 x_i) = 0$$

$$\frac{d}{d\hat{\beta}_1} RSS(\hat{\beta}_0, \hat{\beta}_1) \quad (2)$$

$$= -2 \sum_{i=1}^n x_i (y_i - \hat{\beta}_0 - \hat{\beta}_1 x_i) = 0$$

$$(1) \rightarrow \sum_{i=1}^n y_i - n \hat{\beta}_0 - \hat{\beta}_1 \sum_{i=1}^n x_i = 0$$

$\cancel{n} \div n$

$$\left(\frac{\sum_{i=1}^n y_i}{n} \right) - \hat{\beta}_0 - \hat{\beta}_1 \left(\frac{\sum_{i=1}^n x_i}{n} \right) = 0$$

$$\bar{y} - \hat{\beta}_0 - \hat{\beta}_1 \bar{x} = 0$$

$$\hat{y} = \hat{\beta}_0 + \hat{\beta}_1 \bar{x}$$

(3)
(equivalent
to (1))

$$(2) \Rightarrow \sum_{i=1}^n x_i y_i - \hat{\beta}_0 \sum_{i=1}^n x_i - \hat{\beta}_1 \sum_{i=1}^n x_i^2 = 0$$

$$\div n \Rightarrow \frac{\sum x_i y_i}{n} - \underbrace{\hat{\beta}_0 \bar{x}} - \hat{\beta}_1 \frac{\sum x_i^2}{n} = 0$$

use (3)

$$\Rightarrow \frac{\sum x_i y_i}{n} - (\bar{y} - \hat{\beta}_1 \bar{x}) \bar{x} - \hat{\beta}_1 \frac{\sum x_i^2}{n} = 0$$

$$\frac{\sum x_i y_i}{n} - \bar{x} \bar{y} = \hat{\beta}_1 \left(\frac{\sum x_i^2}{n} - \bar{x}^2 \right)$$

Two relations (exercise)

$$\frac{\sum x_i y_i}{n} - \bar{x} \bar{y} = \frac{1}{n} \sum_{i=1}^n (x_i - \bar{x})(y_i - \bar{y}) \leftarrow$$

$$\frac{\sum x_i^2}{n} - \bar{x}^2 = \frac{1}{n} \sum_{i=1}^n (x_i - \bar{x})^2$$

$$\hat{\beta}_1 = \frac{\sum_{i=1}^n (x_i - \bar{x})(y_i - \bar{y})}{\sum_{i=1}^n (x_i - \bar{x})^2}$$

$$(3) \rightarrow \hat{\beta}_0 = \bar{y} - \hat{\beta}_1 \bar{x}$$

$$\Rightarrow \hat{\beta}_0 = \bar{y} - \bar{x} \frac{\sum_{i=1}^n (x_i - \bar{x})(y_i - \bar{y})}{\sum_{i=1}^n (x_i - \bar{x})^2}$$

Question: How good are the above estimates $\hat{\beta}_0, \hat{\beta}_1$?

To evaluate the performance of our estimates of

$\hat{\beta}_0, \hat{\beta}_1$, let's assume that the data is generated according to the following probabilistic model:

$$y = \underbrace{\beta_0}_{\text{true}} + \underbrace{\beta_1}_{\sim} x + \varepsilon$$

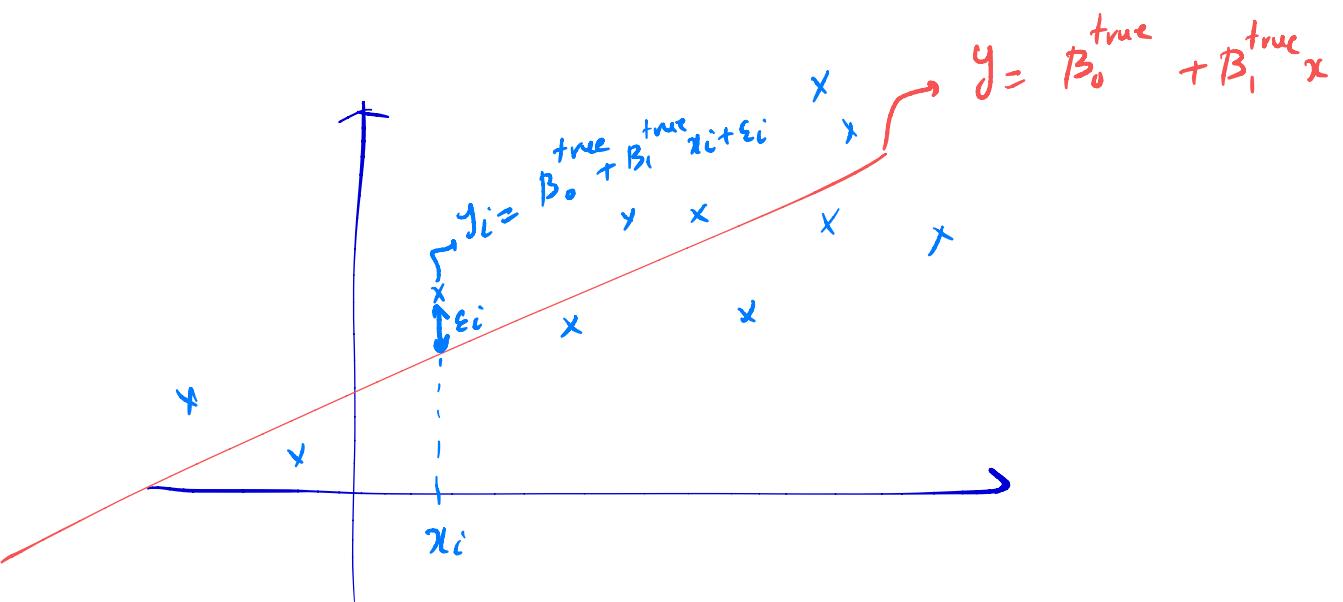
noise
 zero-mean
 independent
 $\text{var}[\varepsilon] = 6^2$

Now, assume that we generate data points

$\{(x_i, y_i)\}_{i=1}^n$ from the above model. Let $\hat{\beta}_0, \hat{\beta}_1$ be the solutions of the regression problem.

as $n \rightarrow \infty$

$\hat{\beta}_0 \xrightarrow{\sim} \beta_0^{\text{true}}$	$\hat{\beta}_0$ is an estimator of: β_0^{true}
$\hat{\beta}_1 \xrightarrow{\sim} \beta_1^{\text{true}}$	$\hat{\beta}_1$ is an estimator of: β_1^{true}



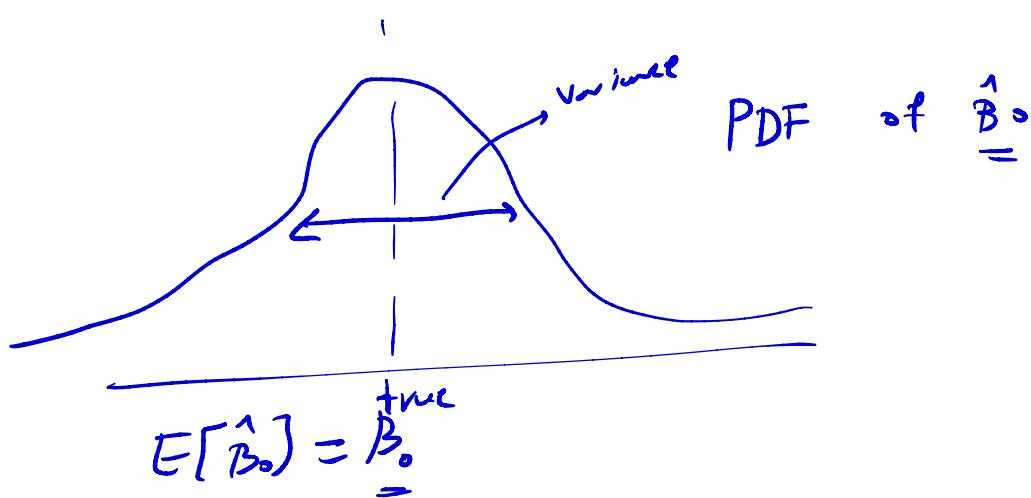
$$\hat{B}_0 \longrightarrow B_0^{\text{true}}$$

$$\hat{B}_1 \longrightarrow B_1^{\text{true}}$$

- In order to quantify the performance of the estimators \hat{B}_0, \hat{B}_1 , we need to compute their expectation (Bias) and Variance.
- The expectation and the variance of \hat{B}_0 (and \hat{B}_1) will tell us how close \hat{B}_0 (and \hat{B}_1) are w.r.t their true value B_0^{true} (and B_1^{true})
- Note that, because of the fact that the noises ϵ_i 's are random variables, our

data set $\{(x_i, y_i)\}$ is also random.

This means that $\hat{\beta}_0, \hat{\beta}_1$, which depend on the data, are random variables as well.



We'll now compute the expectation of $\hat{\beta}_0, \hat{\beta}_1$ and show that they are unbiased estimators of β_0^{true} and β_1^{true} , respectively.

$$E[\hat{\beta}_1] \stackrel{?}{=} \beta_1^{\text{true}}$$

↓

$$E[\hat{\beta}_1] = E\left[\frac{\sum_{i=1}^n (x_i - \bar{x})(y_i - \bar{y})}{\sum_{i=1}^n (x_i - \bar{x})^2} \right]$$

note:

$$y_i = \beta_0^{\text{true}} + \beta_1^{\text{true}} x_i + \varepsilon_i$$

zero-mean
independent
of everything else

Hence:

$$y_i = \beta_0^{\text{true}} + \beta_1^{\text{true}} x_i + \varepsilon_i$$

\Downarrow

$$\bar{y} = \beta_0^{\text{true}} + \beta_1^{\text{true}} \bar{x} + \bar{\varepsilon}$$
(1)

$(\bar{\varepsilon} = \frac{\varepsilon_1 + \dots + \varepsilon_n}{n})$

$$E[\hat{\beta}_1] = E\left[\frac{\sum (x_i - \bar{x})(y_i - \bar{y})}{\sum (x_i - \bar{x})^2} \right]$$

(1)

$$= E\left[\frac{\sum_i (x_i - \bar{x}) (\beta_1^{\text{true}}(x_i - \bar{x}) + \varepsilon_i - \bar{\varepsilon})}{\sum_i (x_i - \bar{x})^2} \right]$$

$$\begin{aligned}
 & (y_i - \bar{y}) \\
 & \downarrow \\
 & (B_0^{\text{true}} + B_1^{\text{true}} x_i + \varepsilon_i) - (B_0^{\text{true}} + B_1^{\text{true}} \bar{x} + \bar{\varepsilon}) \\
 & = B_1^{\text{true}} (x_i - \bar{x}) + \varepsilon_i - \bar{\varepsilon}
 \end{aligned}$$

$$E[\hat{B}_1] = E \left[\frac{B_1^{\text{true}} \sum_i (x_i - \bar{x})^2 + \sum (\varepsilon_i - \bar{\varepsilon})(x_i - \bar{x})}{\sum_i (x_i - \bar{x})^2} \right]$$

$$= B_1^{\text{true}} + E \left[\sum_{i=1}^n \frac{(\varepsilon_i - \bar{\varepsilon})(x_i - \bar{x})}{\sum_i (x_i - \bar{x})^2} \right]$$

noises
 were independent
 of everything
 else

$$= B_1^{\text{true}} + \sum_{i=1}^n E \left[\frac{(\varepsilon_i - \bar{\varepsilon})(x_i - \bar{x})}{\sum_i (x_i - \bar{x})^2} \right]$$

$$= B_1^{\text{true}} + \sum_{i=1}^n \underbrace{E[\varepsilon_i - \bar{\varepsilon}]}_{E[\varepsilon_i] - E[\bar{\varepsilon}]} \underbrace{E \left[\frac{(x_i - \bar{x})}{\sum_i (x_i - \bar{x})^2} \right]}_{0}$$

$$\Rightarrow E[\hat{\beta}_1] = \beta_1^{\text{true}}$$

In general:

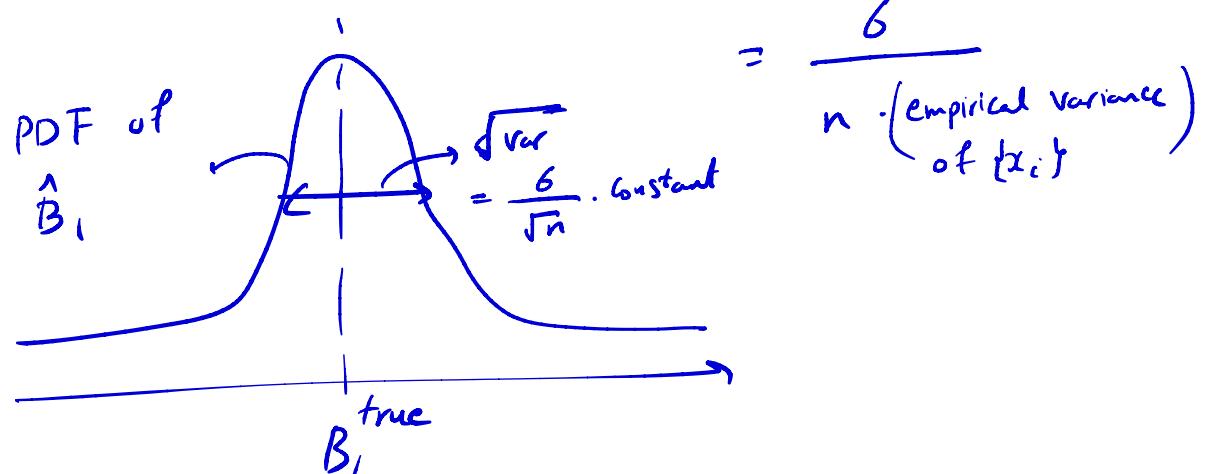
$$(1) \quad E[\hat{\beta}_1] = \beta_1^{\text{true}}$$

$$(2) \quad E[\hat{\beta}_0] = \beta_0^{\text{true}}$$

$$(3) \quad \text{Var}(\hat{\beta}_0) = E[(\hat{\beta}_0 - \beta_0^{\text{true}})^2]$$

$$(\text{Var}(\varepsilon_i) = \sigma^2) \quad = \sigma^2 \left(\frac{1}{n} + \frac{\bar{x}^2}{\sum_{i=1}^n (x_i - \bar{x})^2} \right)$$

$$(4) \quad \text{Var}(\hat{\beta}_1) = \frac{\sigma^2}{\sum_{i=1}^n (x_i - \bar{x})^2} = \frac{\sigma^2}{n \underbrace{\left(\frac{1}{n} \sum_{i=1}^n (x_i - \bar{x})^2 \right)}_{\text{empirical variance of } \{x_i\}}}$$



So far we've only been talking about the 1-dim case ($p=1$). Let's now consider the general case:

$$\{ (x_i, y_i) \}_{i=1}^n \quad \begin{array}{l} x_i \in \mathbb{R}^p, y_i \in \mathbb{R} \\ \text{---} \\ x_i = \begin{pmatrix} x_{i1} \\ \vdots \\ x_{ip} \end{pmatrix} \in \mathbb{R}^p, y_i \in \mathbb{R} \end{array}$$

simplest parametric family (linear):

$$y_i \approx \hat{\beta}_0 + \hat{\beta}_1 x_{i1} + \hat{\beta}_2 x_{i2} + \cdots + \hat{\beta}_p x_{ip}$$

$\underbrace{\qquad\qquad\qquad}_{f(x_i, \vec{\beta})}$