

6 Oct 2019

$\in \mathbb{N}$

## Lecture 2

Data  $\equiv$  [EDA]

[counts in discrete space]  
[intervals:]

I.  $\left\{ \begin{array}{l} \text{categorical} \\ \text{quantitative} \end{array} \right\} \rightarrow \text{classes}$   
 $\left\{ \begin{array}{l} \text{categorical} \\ \text{quantitative} \end{array} \right\} \rightarrow \text{numerical}$   
 $x, y \in \mathbb{R}$   
 $x \in \mathbb{N}$   
 $\vec{x} \in \mathbb{R}^d$   
 $X \in \mathbb{R}^{m \times n}$

$\left\{ \begin{array}{l} \text{continuous} \\ \text{discrete} \\ \text{ordinals} \end{array} \right\} \rightarrow \mathbb{R}$   
 $\rightarrow \mathbb{N}$

II.  $\left[ \begin{array}{l} \text{univariate} \\ \updownarrow \\ \text{multivariate} \end{array} \right] \mathbb{R}^d$

III. Predictors  $\Rightarrow$  Outcome

dependency: [independent]  $\Rightarrow$  [dependent]

[control input]  $\Rightarrow$  [observable output]

Both can have randomness  
but

$x + \epsilon_{in} \rightsquigarrow f(x + \epsilon_{in}) + \delta_{obs}$

Learning: from input  $\rightarrow$  output

# Brief summary for supervised learning

Given data  $\longrightarrow$  EDA (clean-up)  
 $\downarrow$   
decide on variables  
+ task  
 $\underbrace{\hspace{10em}}$

$\longrightarrow$  Declare predictor space  $X$  ; outcome space  $Y$

$\longrightarrow$   $N$  <sup>many</sup> samples to be used for learning  $\left\{ \begin{array}{l} (x_i, y_i) \quad i=1, \dots, N \\ x_i \in X \quad y_i \in Y \end{array} \right.$

In Supervised learning ~~means that~~  
this is the training set.

Keep some samples for validation

$\nwarrow$  For a given task with training set  $\{ (x_i, y_i) \}_{i=1}^N$

$\Rightarrow$  Find function

$$f: X \rightarrow Y$$

and define

loss function  $L(f(x), y)$

such that

(i) in-sample loss (error) is minimised

$$\mathbb{E} [L(f(\{x_i\}), \{y_i\})] \ll$$

and  $\{(x_i, y_i)\}_{i \in \text{training}}$

(ii) expected out-of-sample loss is also small

$$\mathbb{E} [L(f(\{x_k\}), \{y_k\})] \ll$$

$\{(x_k, y_k)\}_{k \in \text{test}}$

(validation)

$\Downarrow$   
we expect then that

$f(x_{\text{unknown}})$  will be a  
good predictor for  $y_{\text{unknown}}$

# Linear regression

Data:                      Inputs  
                                    (predictors)  
                                     $X$

                                    Output  
                                    (outcome)  
                                     $y$

Quantitative

Samples :

$i = 1, \dots, N$

samples  
(observations)

$x_1^{(1)}$	$x_2^{(1)}$	...	$x_p^{(1)}$	$y^{(1)}$
$x_1^{(2)}$	$x_2^{(2)}$	...	$x_p^{(2)}$	$y^{(2)}$
$\vdots$	$\vdots$		$\vdots$	
$x_1^{(N)}$	$x_2^{(N)}$	...	$x_p^{(N)}$	$y^{(N)}$
Inputs : $\{ \vec{x}^{(i)} \}_{i=1}^N$				output $\{ y^{(i)} \}_{i=1}^N$

Assume a linear relationship:  $f(\vec{x}) = \beta_0 + \beta_1 x_1 + \beta_2 x_2 + \dots + \beta_p x_p$

$\vec{x} \in \mathbb{R}^p$   
 $y \in \mathbb{R}$

$$\hat{y} = f_{LR}(\vec{x}; \vec{\theta})$$

$\vec{\theta} = (\beta_0, \dots, \beta_p)$   
↑ (Hyp) Parameters

Need to find  $\vec{\theta}$

Defines the model

$f: X \rightarrow Y$   
 $f: \mathbb{R}^p \rightarrow \mathbb{R}$

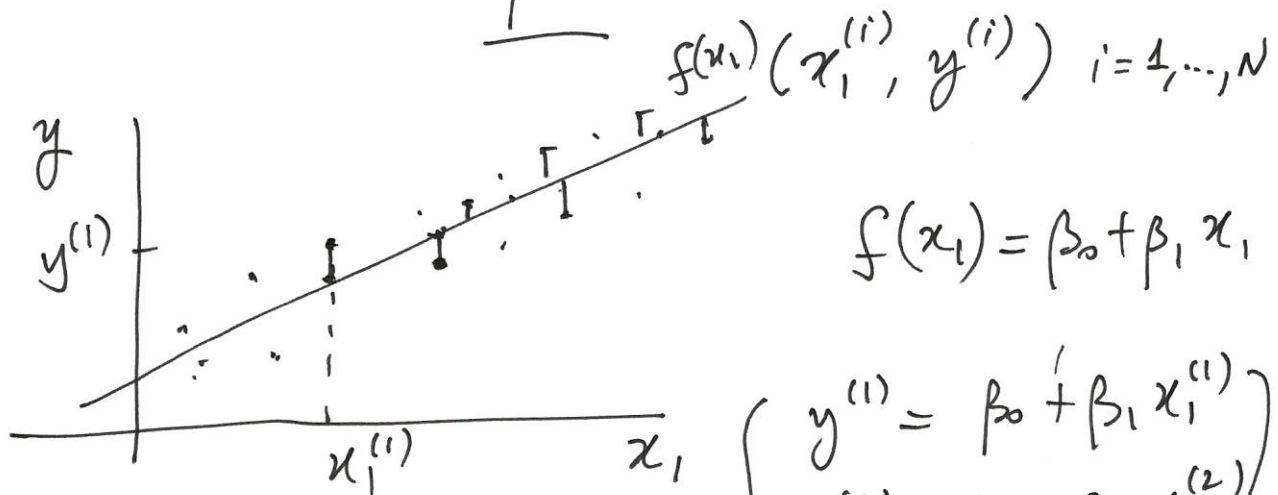


Simplest case (so that we can draw).

Assume

$$\underline{p=1}$$

$$\underline{N \gg 2}$$



Mean Square error :

MSE

$$MSE = \frac{1}{N} \sum_{i=1}^N \left[ y^{(i)} - f_{LR}(x_1^{(i)}; \beta_0, \beta_1) \right]^2 = \frac{1}{N} \sum_{i=1}^N e^{(i)2}$$

$L(y, f_{LR}(x_1))$

Define :

$$\vec{y}^T = (y^{(1)}, \dots, y^{(N)}) \in \mathbb{R}^N$$

$$\vec{x}_1^T = (x_1^{(1)}, \dots, x_1^{(N)}) \in \mathbb{R}^N$$

$$\vec{1}^T = (1, \dots, 1) \in \mathbb{R}^N$$

$$f_{LR}(\vec{x}) = \beta_0 \vec{1}_{N \times 1} + \beta_1 \vec{x}_1_{N \times 1} = \begin{pmatrix} \vec{1} & \vec{x}_1 \end{pmatrix}_{N \times 2} \begin{pmatrix} \beta_0 \\ \beta_1 \end{pmatrix}_{2 \times 1}$$

$\equiv$   
 $\times$

$$f_{LR}(\vec{x}) = X \cdot \vec{\beta} \quad \text{Reminder} \quad (\vec{\beta} = \vec{\theta})$$

$$\vec{e} = \vec{y} - X\vec{\beta}$$

$$L[f_{LR}(\vec{x}_i), y] = \frac{1}{N} \vec{e}^T \cdot \vec{e} = \frac{\|\vec{e}\|^2}{N}$$

Minimize  $L$  in the space of parameter  $\vec{\beta}$

$$L(\vec{\beta}) = L(\beta_0, \beta_1) = \frac{1}{N} [(\vec{y} - X\vec{\beta})^T (\vec{y} - X\vec{\beta})]$$

$$(1) \quad \left. \frac{dL}{d\vec{\beta}} \right|_{\vec{\beta}^*} = \nabla_{\vec{\beta}} L \Big|_{\vec{\beta}^*} = 0$$

$$(2) \text{ Check that } H(\vec{\beta}^*)_{ij} = \left( \frac{\partial^2 L}{\partial \beta_i \partial \beta_j} \right) \quad H > 0$$

$$L(\vec{\beta}) = \frac{1}{N} [\vec{y}^T \vec{y} - \vec{y}^T X \vec{\beta} - \vec{\beta}^T X^T \vec{y} + \vec{\beta}^T X^T X \vec{\beta}]$$

$$\nabla_{\vec{\beta}} L(\vec{y}^T X \vec{\beta})$$

Aside

$$\nabla_{\vec{\beta}} (\vec{\alpha}^T \vec{\beta}) = \nabla_{\vec{\beta}} (\alpha_1 \beta_0 + \alpha_2 \beta_1) = \begin{pmatrix} \alpha_1 \\ \alpha_2 \end{pmatrix} \equiv \vec{\alpha}$$

$$\nabla_{\vec{\beta}} (\vec{\beta}^T \vec{\alpha}) = \vec{\alpha}$$

$$\nabla_{\vec{\beta}} (\vec{\beta}^T A \vec{\beta}) = A \vec{\beta} + A^T \vec{\beta} = (A + A^T) \vec{\beta}$$

$$\begin{aligned} \nabla_{\vec{\beta}} L(\vec{\beta}) &= \frac{1}{N} [-X^T \vec{y} - X^T \vec{y} + (X^T X + (X^T X)^T) \vec{\beta}] \\ &= \frac{-2}{N} [\vec{y}^T X \vec{\beta} - X^T X \vec{\beta}] \quad \text{Normal equations.} \end{aligned}$$

$$\nabla_{\vec{\beta}} L|_{\vec{\beta}^*} = 0$$

$$\underline{X^T \vec{y} = (X^T X) \vec{\beta}^*}$$

$$X = \begin{bmatrix} \vec{1} & \vec{x}_1 \end{bmatrix}_{N \times 2} \text{ invertible}$$

$$\underline{\vec{\beta}^* = (X^T X)^{-1} X^T \vec{y}}$$