Big Data Analysis (MA60306)

Bibhas Adhikari

Spring 2022-23, IIT Kharagpur

Lecture 4 January 11, 2023

Bibhas Adhikari (Spring 2022-23, IIT Kharag

Big Data Analysis

Lecture 4 January 11, 2023 1 / 11

3

Plots for two quantitative variables

Example of scatterplot of height and weight data of the nutri data

#### Plots for two quantitative variables

Example of scatterplot of height and weight data of the *nutri* data Example of scatterplot of the *Advertising* data

Statistical Learning Mathematical analysis of the data - interpret the model and quantify the uncertainty in the data

3

Statistical Learning Mathematical analysis of the data - interpret the model and quantify the uncertainty in the data

Machine Learning - the emphasis is on making predictions using large-scale data

3

Statistical Learning Mathematical analysis of the data - interpret the model and quantify the uncertainty in the data

Machine Learning - the emphasis is on making predictions using large-scale data

Modeling data

▷ accurately predict some future quantity of interest, given some observed data - prediction - exp. direct-marketing campaign

Statistical Learning Mathematical analysis of the data - interpret the model and quantify the uncertainty in the data

Machine Learning - the emphasis is on making predictions using large-scale data

Modeling data

- accurately predict some future quantity of interest, given some observed data - prediction - exp. direct-marketing campaign
- $\triangleright\,$  discover patterns in the data inference exp. advertising

Statistical Learning Mathematical analysis of the data - interpret the model and quantify the uncertainty in the data

Machine Learning - the emphasis is on making predictions using large-scale data

Modeling data

- accurately predict some future quantity of interest, given some observed data - prediction - exp. direct-marketing campaign
- $\triangleright\,$  discover patterns in the data inference exp. advertising
- ▷ both prediction and inference exp. real estate data

< ロ > < 同 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ >

Statistical Learning Mathematical analysis of the data - interpret the model and quantify the uncertainty in the data

Machine Learning - the emphasis is on making predictions using large-scale data

Modeling data

- accurately predict some future quantity of interest, given some observed data - prediction - exp. direct-marketing campaign
- ▷ discover patterns in the data inference exp. advertising
- $\triangleright\,$  both prediction and inference exp. real estate data

Tools:

 Functional approximation - how one data variable depends on another data variable

イロト 不得 トイラト イラト 一日

Statistical Learning Mathematical analysis of the data - interpret the model and quantify the uncertainty in the data

Machine Learning - the emphasis is on making predictions using large-scale data

Modeling data

- accurately predict some future quantity of interest, given some observed data - prediction - exp. direct-marketing campaign
- ▷ discover patterns in the data inference exp. advertising
- $\triangleright\,$  both prediction and inference exp. real estate data

Tools:

- Functional approximation how one data variable depends on another data variable
- > Optimization best possible model in a class of models

Prediction function g: Let x be a feature/input vector. Then one of the fundamental problems in machine learning is to predict an output response variable y. The prediction for y is based on a function g(x) such that g encompasses all the information about the relationship between the variables x and y. - exp. digitized signature

Prediction function g: Let x be a feature/input vector. Then one of the fundamental problems in machine learning is to predict an output response variable y. The prediction for y is based on a function g(x) such that g encompasses all the information about the relationship between the variables x and y. - exp. digitized signature

Regression In regression problems, the response variable y can take any real value

< ロ > < 同 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ >

Prediction function g: Let x be a feature/input vector. Then one of the fundamental problems in machine learning is to predict an output response variable y. The prediction for y is based on a function g(x) such that g encompasses all the information about the relationship between the variables x and y. - exp. digitized signature

Regression In regression problems, the response variable y can take any real value

Classification If y lies in a finite set then predicting  $y \in \{0, ..., k-1\}$  is same as classifying the input into one of the k categories

イロト 不得 トイラト イラト 一日

Prediction function g: Let x be a feature/input vector. Then one of the fundamental problems in machine learning is to predict an output response variable y. The prediction for y is based on a function g(x) such that g encompasses all the information about the relationship between the variables x and y. - exp. digitized signature

Regression In regression problems, the response variable y can take any real value

Classification If y lies in a finite set then predicting  $y \in \{0, ..., k-1\}$  is same as classifying the input into one of the k categories

Loss function Let  $\hat{y}$  denote a prediction corresponding to a given response. Then we do we measure the accuracy of the prediction? We write  $Loss(y, \hat{y})$  to denote the measure

▷ regression: 
$$Loss(y, \hat{y}) = (y - \hat{y})^2$$

 $\triangleright \text{ classification: } Loss(y, \widehat{y}) = \mathbb{1}\{y \neq \widehat{y}\}$ 

$$\triangleright$$
 in general,  $Loss(y, \widehat{y}) = ||y - \widehat{y}||^2$ 

## Supervised learning technique

The word 'regression' - In 1889, Francis Galton observed that the heights of the adult offsprings have more average heights/intelligence than their parents - used the term 'regression' to indicate "return to mediocrity"

## Supervised learning technique

The word 'regression' - In 1889, Francis Galton observed that the heights of the adult offsprings have more average heights/intelligence than their parents - used the term 'regression' to indicate "return to mediocrity"

Let  $\mathbf{x} = [x_1, \dots, x_p]^T$  be an input with p features and we need to predict a quantitative response/output y via a function  $g(\mathbf{x})$ ,  $x_i$  could be discrete or continuous

## Supervised learning technique

The word 'regression' - In 1889, Francis Galton observed that the heights of the adult offsprings have more average heights/intelligence than their parents - used the term 'regression' to indicate "return to mediocrity"

Let  $\mathbf{x} = [x_1, \dots, x_p]^T$  be an input with p features and we need to predict a quantitative response/output y via a function  $g(\mathbf{x})$ ,  $x_i$  could be discrete or continuous

#### Example

- predict birth weight of a baby from the weight of the mother, her socio-economic status, her smoking habits (sometimes known as explanatory variables)
- Consider the advertising data for advertising budgets for the product in each of 200 markets for three different media: TV, radio, and newspaper. Can we develop a model to predict sales based on the three media budgets?

イロト 不得 トイヨト イヨト 二日

### Defining it as a statistical problem

 $\triangleright\,$  It is unlikely that the prediction function g will make accurate prediction

э

### Defining it as a statistical problem

- $\triangleright\,$  It is unlikely that the prediction function g will make accurate prediction
- $\triangleright$  in reality, even for the same input **x**, the output y may be different

э

### Defining it as a statistical problem

- $\triangleright\,$  It is unlikely that the prediction function g will make accurate prediction
- $\triangleright$  in reality, even for the same input **x**, the output y may be different
- ▷ we adopt a probabilistic approach: we assume that each pair  $(\mathbf{x}, y)$  is an instance of a random variable  $(\mathbf{X}, Y)$  that has the joint pdf  $f(\mathbf{x}, y)$

### Defining it as a statistical problem

- $\triangleright\,$  It is unlikely that the prediction function g will make accurate prediction
- $\triangleright$  in reality, even for the same input **x**, the output y may be different
- ▷ we adopt a probabilistic approach: we assume that each pair  $(\mathbf{x}, y)$  is an instance of a random variable  $(\mathbf{X}, Y)$  that has the joint pdf  $f(\mathbf{x}, y)$
- ▷ Then the performance of the prediction can be measured via the expected loss, known as *risk*:

$$I(g) = \mathbb{E}[\text{Loss}(Y, g(\mathbf{X}))]$$

### Defining it as a statistical problem

- $\triangleright\,$  It is unlikely that the prediction function g will make accurate prediction
- $\triangleright\,$  in reality, even for the same input  ${\bf x},$  the output y may be different
- ▷ we adopt a probabilistic approach: we assume that each pair  $(\mathbf{x}, y)$  is an instance of a random variable  $(\mathbf{X}, Y)$  that has the joint pdf  $f(\mathbf{x}, y)$
- ▷ Then the performance of the prediction can be measured via the expected loss, known as *risk*:

$$I(g) = \mathbb{E}[\text{Loss}(Y, g(\mathbf{X}))]$$

▷ In the classification case, the risk equals the probability of incorrect classification:

$$I(g) = \mathbb{P}[Y \neq g(\mathbf{X})]$$

### Defining it as a statistical problem

- $\triangleright\,$  It is unlikely that the prediction function g will make accurate prediction
- $\triangleright$  in reality, even for the same input **x**, the output y may be different
- $\triangleright$  we adopt a probabilistic approach: we assume that each pair  $(\mathbf{x}, y)$  is an instance of a random variable  $(\mathbf{X}, Y)$  that has the joint pdf  $f(\mathbf{x}, y)$
- ▷ Then the performance of the prediction can be measured via the expected loss, known as *risk*:

$$I(g) = \mathbb{E}[\text{Loss}(Y, g(\mathbf{X}))]$$

▷ In the classification case, the risk equals the probability of incorrect classification:

$$l(g) = \mathbb{P}[Y \neq g(\mathbf{X})]$$

In this case, g is called the classifier

Image: A matrix

Given a rv  $(\mathbf{X}, Y)$  and a loss function, in principle, we determine

$$g^* = \arg\min_g \mathbb{E}[\operatorname{Loss}(Y, g(\mathbf{X}))]$$

which yields the smallest risk  $I^* = I(g^*)$ 

3

Given a rv  $(\mathbf{X}, Y)$  and a loss function, in principle, we determine

$$g^* = rg\min_g \mathbb{E}[\mathsf{Loss}(Y, g(\mathsf{X}))]$$

which yields the smallest risk  $I^* = I(g^*)$ 

Theorem If  $Loss(y, \hat{y}) = (y - \hat{y})^2$  then the optimal prediction function  $g^*$  is equal to the conditional expectation of Y given  $\mathbf{X} = \mathbf{x}$ :

$$g^*(\mathsf{x}) = \mathbb{E}[Y|\mathsf{X} = \mathsf{x}]$$

Given a rv  $(\mathbf{X}, Y)$  and a loss function, in principle, we determine

$$g^* = rg \min_g \mathbb{E}[\mathsf{Loss}(Y, g(\mathsf{X}))]$$

which yields the smallest risk  $I^* = I(g^*)$ 

Theorem If  $Loss(y, \hat{y}) = (y - \hat{y})^2$  then the optimal prediction function  $g^*$  is equal to the conditional expectation of Y given  $\mathbf{X} = \mathbf{x}$ :

$$g^*(\mathsf{x}) = \mathbb{E}[Y|\mathsf{X} = \mathsf{x}]$$

Proof

$$\mathbb{E}[(Y-g(\mathbf{X}))^2] = \mathbb{E}[(Y-g^*(\mathbf{X})+g^*(\mathbf{X})-g(\mathbf{X}))^2]$$

・ 何 ト ・ ヨ ト ・ ヨ ト

Given a rv  $(\mathbf{X}, Y)$  and a loss function, in principle, we determine

$$g^* = rg\min_g \mathbb{E}[\mathsf{Loss}(Y, g(\mathsf{X}))]$$

which yields the smallest risk  $I^* = I(g^*)$ 

Theorem If  $Loss(y, \hat{y}) = (y - \hat{y})^2$  then the optimal prediction function  $g^*$  is equal to the conditional expectation of Y given  $\mathbf{X} = \mathbf{x}$ :

$$g^*(\mathsf{x}) = \mathbb{E}[Y|\mathsf{X} = \mathsf{x}]$$

Proof

$$\begin{split} \mathbb{E}[(Y - g(\mathbf{X}))^2] &= \mathbb{E}[(Y - g^*(\mathbf{X}) + g^*(\mathbf{X}) - g(\mathbf{X}))^2] \\ &= \mathbb{E}[(Y - g^*(\mathbf{X}))^2] + 2\mathbb{E}[(Y - g^*(\mathbf{X}))(g^*(\mathbf{X}) - g(\mathbf{X}))] \\ &+ \mathbb{E}[(g^*(\mathbf{X}) - g(\mathbf{X}))^2] \end{split}$$

・ 同 ト ・ ヨ ト ・ ヨ ト

Given a rv  $(\mathbf{X}, Y)$  and a loss function, in principle, we determine

$$g^* = rg\min_g \mathbb{E}[\mathsf{Loss}(Y, g(\mathsf{X}))]$$

which yields the smallest risk  $I^* = I(g^*)$ 

Theorem If  $Loss(y, \hat{y}) = (y - \hat{y})^2$  then the optimal prediction function  $g^*$  is equal to the conditional expectation of Y given  $\mathbf{X} = \mathbf{x}$ :

$$g^*(\mathsf{x}) = \mathbb{E}[Y|\mathsf{X} = \mathsf{x}]$$

#### Proof

$$\begin{split} \mathbb{E}[(Y - g(\mathbf{X}))^2] &= \mathbb{E}[(Y - g^*(\mathbf{X}) + g^*(\mathbf{X}) - g(\mathbf{X}))^2] \\ &= \mathbb{E}[(Y - g^*(\mathbf{X}))^2] + 2\mathbb{E}[(Y - g^*(\mathbf{X}))(g^*(\mathbf{X}) - g(\mathbf{X}))] \\ &+ \mathbb{E}[(g^*(\mathbf{X}) - g(\mathbf{X}))^2] \\ &\geq \mathbb{E}[(Y - g^*(\mathbf{X}))^2] + 2\mathbb{E}[(Y - g^*(\mathbf{X}))(g^*(\mathbf{X}) - g(\mathbf{X}))] \end{split}$$

・ 何 ト ・ ヨ ト ・ ヨ ト

Given a rv  $(\mathbf{X}, Y)$  and a loss function, in principle, we determine

$$g^* = rg\min_g \mathbb{E}[\mathsf{Loss}(Y, g(\mathsf{X}))]$$

which yields the smallest risk  $I^* = I(g^*)$ 

Theorem If  $Loss(y, \hat{y}) = (y - \hat{y})^2$  then the optimal prediction function  $g^*$  is equal to the conditional expectation of Y given  $\mathbf{X} = \mathbf{x}$ :

$$g^*(\mathsf{x}) = \mathbb{E}[Y|\mathsf{X} = \mathsf{x}]$$

#### Proof

$$\begin{split} & \mathbb{E}[(Y - g(\mathbf{X}))^2] = \mathbb{E}[(Y - g^*(\mathbf{X}) + g^*(\mathbf{X}) - g(\mathbf{X}))^2] \\ &= \mathbb{E}[(Y - g^*(\mathbf{X}))^2] + 2\mathbb{E}[(Y - g^*(\mathbf{X}))(g^*(\mathbf{X}) - g(\mathbf{X}))] \\ &+ \mathbb{E}[(g^*(\mathbf{X}) - g(\mathbf{X}))^2] \\ &\geq \mathbb{E}[(Y - g^*(\mathbf{X}))^2] + 2\mathbb{E}[(Y - g^*(\mathbf{X}))(g^*(\mathbf{X}) - g(\mathbf{X}))] \\ &= \mathbb{E}[(Y - g^*(\mathbf{X}))^2] + 2\mathbb{E}\{(g^*(\mathbf{X}) - g(\mathbf{X}))\mathbb{E}[Y - g^*(\mathbf{X})|\mathbf{X}]\} \end{split}$$

7/11

Now

$$\mathbb{E}[Y - g^*(\mathbf{X})|\mathbf{X}] = 0$$

and hence

$$\mathbb{E}[(Y - g(\mathbf{X}))^2] \geq \mathbb{E}[(Y - g^*(\mathbf{X}))^2]$$

which implies  $g^*$  yields the smallest squared-error.

3

Now

$$\mathbb{E}[Y - g^*(\mathbf{X}) | \mathbf{X}] = 0$$

and hence

$$\mathbb{E}[(Y - g(\mathbf{X}))^2] \geq \mathbb{E}[(Y - g^*(\mathbf{X}))^2]$$

which implies  $g^*$  yields the smallest squared-error.

#### Consequence

 $\triangleright$  The conditional **X** = **x**, the random response *Y* can be written as

$$Y = g^*(\mathbf{x}) + \epsilon(\mathbf{x})$$

where  $\epsilon(\mathbf{x})$  can be thought of as a random deviation of the response from its conditional mean at  $\mathbf{x}$ .

Now

$$\mathbb{E}[Y - g^*(\mathbf{X}) | \mathbf{X}] = 0$$

and hence

$$\mathbb{E}[(Y - g(\mathbf{X}))^2] \geq \mathbb{E}[(Y - g^*(\mathbf{X}))^2]$$

which implies  $g^*$  yields the smallest squared-error.

#### Consequence

 $\triangleright$  The conditional **X** = **x**, the random response Y can be written as

$$Y = g^*(\mathbf{x}) + \epsilon(\mathbf{x})$$

where  $\epsilon(\mathbf{x})$  can be thought of as a random deviation of the response from its conditional mean at  $\mathbf{x}$ .

 $\triangleright \ \mathbb{E}[\epsilon(\mathbf{x})] = \mathbf{0}$ 

Now

$$\mathbb{E}[Y-g^*(\mathbf{X})|\mathbf{X}]=0$$

and hence

$$\mathbb{E}[(Y - g(\mathbf{X}))^2] \geq \mathbb{E}[(Y - g^*(\mathbf{X}))^2]$$

which implies  $g^*$  yields the smallest squared-error.

#### Consequence

 $\triangleright$  The conditional **X** = **x**, the random response Y can be written as

$$Y = g^*(\mathbf{x}) + \epsilon(\mathbf{x})$$

where  $\epsilon(\mathbf{x})$  can be thought of as a random deviation of the response from its conditional mean at  $\mathbf{x}$ .

$$\label{eq:expansion} \begin{array}{l} \triangleright \ \mathbb{E}[\epsilon(\mathbf{x})] = \mathbf{0} \\ \ \triangleright \ \mathbb{V}\mathrm{ar}[\epsilon(\mathbf{x})] = \nu^2(\mathbf{x}) \mbox{ for some function } \nu(\mathbf{x}) \end{array}$$

### Learning

 $\triangleright$  The optimal prediction function  $g^*$  depends on the joint distribution of  $(\mathbf{X}, Y)$ , which is not available in practice

э

### Learning

- $\triangleright$  The optimal prediction function  $g^*$  depends on the joint distribution of  $(\mathbf{X}, Y)$ , which is not available in practice
- $\triangleright$  available : finite number of (usually) independent realizations from the joint density  $f(\mathbf{x}, y)$

### Learning

- $\triangleright$  The optimal prediction function  $g^*$  depends on the joint distribution of  $(\mathbf{X}, Y)$ , which is not available in practice
- $\triangleright$  available : finite number of (usually) independent realizations from the joint density  $f(\mathbf{x}, y)$
- ▷ Let  $\mathcal{T} = \{(\mathbf{X}_1, Y_1), \dots, (\mathbf{X}_n, Y_n)\}$  be a sample, and call it the *training* set with *n* examples

A B A A B A

### Learning

- $\triangleright$  The optimal prediction function  $g^*$  depends on the joint distribution of  $(\mathbf{X}, Y)$ , which is not available in practice
- $\triangleright$  available : finite number of (usually) independent realizations from the joint density  $f(\mathbf{x}, y)$
- ▷ Let  $\mathcal{T} = \{(\mathbf{X}_1, Y_1), \dots, (\mathbf{X}_n, Y_n)\}$  be a sample, and call it the *training* set with *n* examples
- $\triangleright$  a given sample point:  $\{(\mathbf{x}_1, y_1), \dots, (\mathbf{x}_n, y_n)\}$

### Learning

- $\triangleright$  The optimal prediction function  $g^*$  depends on the joint distribution of  $(\mathbf{X}, Y)$ , which is not available in practice
- $\triangleright$  available : finite number of (usually) independent realizations from the joint density  $f(\mathbf{x}, y)$
- ▷ Let  $\mathcal{T} = \{(\mathbf{X}_1, Y_1), \dots, (\mathbf{X}_n, Y_n)\}$  be a sample, and call it the *training* set with *n* examples
- $\triangleright$  a given sample point:  $\{(\mathbf{x}_1, y_1), \dots, (\mathbf{x}_n, y_n)\}$
- $\triangleright$  Goal: 'learn' the unknown function  $g^*$  using the n examples from the training set  ${\cal T}$

- 3

### Learning

- $\triangleright$  The optimal prediction function  $g^*$  depends on the joint distribution of  $(\mathbf{X}, Y)$ , which is not available in practice
- $\triangleright$  available : finite number of (usually) independent realizations from the joint density  $f(\mathbf{x}, y)$
- ▷ Let  $\mathcal{T} = \{(\mathbf{X}_1, Y_1), \dots, (\mathbf{X}_n, Y_n)\}$  be a sample, and call it the *training* set with *n* examples
- $\triangleright$  a given sample point:  $\{(\mathbf{x}_1, y_1), \dots, (\mathbf{x}_n, y_n)\}$
- $\triangleright$  Goal: 'learn' the unknown function  $g^*$  using the n examples from the training set  ${\cal T}$
- $\triangleright\,$  Denote  $g_{\mathcal{T}}$  as an approximation for  $g^*$  that can be constructed from  $\mathcal{T}$

イロト 不得 トイラト イラト 一日

### Learning

- $\triangleright$  The optimal prediction function  $g^*$  depends on the joint distribution of  $(\mathbf{X}, Y)$ , which is not available in practice
- $\triangleright$  available : finite number of (usually) independent realizations from the joint density  $f(\mathbf{x}, y)$
- ▷ Let  $\mathcal{T} = \{(\mathbf{X}_1, Y_1), \dots, (\mathbf{X}_n, Y_n)\}$  be a sample, and call it the *training* set with *n* examples
- $\triangleright$  a given sample point:  $\{(\mathbf{x}_1, y_1), \dots, (\mathbf{x}_n, y_n)\}$
- $\triangleright$  Goal: 'learn' the unknown function  $g^*$  using the n examples from the training set  ${\mathcal T}$
- $\triangleright$  Denote  $g_T$  as an approximation for  $g^*$  that can be constructed from TObviously,  $g_T$  is a random function and a particular outcome is  $g_\tau$

Linear regression model The response Y depends on a *d*-dimensional explanatory vector  $\mathbf{x} = [x_1, \dots, x_d]^T$  via the linear relationship

$$Y = \beta_0 + \beta_1 x_1 + \ldots + \beta_d x_d + \epsilon = \boldsymbol{\beta}^T \mathbf{x} + \epsilon = \mathbf{x}^T \boldsymbol{\beta} + \epsilon$$

3

Linear regression model The response Y depends on a *d*-dimensional explanatory vector  $\mathbf{x} = [x_1, \dots, x_d]^T$  via the linear relationship

$$Y = \beta_0 + \beta_1 x_1 + \ldots + \beta_d x_d + \epsilon = \boldsymbol{\beta}^T \mathbf{x} + \epsilon = \mathbf{x}^T \boldsymbol{\beta} + \epsilon := g(\mathbf{x}) + \epsilon,$$

where  $\mathbb{E}[\epsilon] = 0$  and  $\mathbb{V}ar[\epsilon] = \sigma^2$ .

- 3

Linear regression model The response Y depends on a *d*-dimensional explanatory vector  $\mathbf{x} = [x_1, \dots, x_d]^T$  via the linear relationship

$$Y = \beta_0 + \beta_1 x_1 + \ldots + \beta_d x_d + \epsilon = \boldsymbol{\beta}^T \mathbf{x} + \epsilon = \mathbf{x}^T \boldsymbol{\beta} + \epsilon := g(\mathbf{x}) + \epsilon,$$

where  $\mathbb{E}[\epsilon] = 0$  and  $\mathbb{V}ar[\epsilon] = \sigma^2$ .

The model for a training set  $\mathcal{T} = \{(\mathbf{x}_1, Y_1), \dots, (\mathbf{x}_n, Y_n)\}$  is given by

$$\mathbf{Y} = \mathbf{X}oldsymbol{eta} + oldsymbol{\epsilon}$$

where  $\boldsymbol{\epsilon} = [\epsilon_1, \dots, \epsilon_n]^T$  is a vector of iid copies of  $\epsilon$ , and **X** is the *model* matrix or regression matrix given by

$$\mathbf{X} = \begin{bmatrix} 1 & x_{11} & x_{12} & \dots & x_{1d} \\ 1 & x_{21} & x_{22} & \dots & x_{2d} \\ \vdots & \vdots & \vdots & \vdots & \vdots \\ 1 & x_{n1} & x_{n2} & \dots & x_{nd} \end{bmatrix} = \begin{bmatrix} 1 & \mathbf{x}_1^T \\ 1 & \mathbf{x}_2^T \\ \vdots & \vdots \\ 1 & \mathbf{x}_n^T \end{bmatrix}$$

▷ In most cases  $x_{ij}$ s are generally chosen so that the columns of **X** are linearly independent i.e rank of **X** is d + 1

3

(日) (四) (日) (日) (日)

- $\triangleright$  In most cases  $x_{ij}$ s are generally chosen so that the columns of **X** are linearly independent i.e rank of **X** is d + 1
- ▷ However, in some experimental design situations  $x_{ij} \in \{0, 1\}$  and columns of **X** are linearly dependent, and **X** is called the *design matrix*

A B A A B A

- ▷ In most cases  $x_{ij}$ s are generally chosen so that the columns of **X** are linearly independent i.e rank of **X** is d + 1
- ▷ However, in some experimental design situations  $x_{ij} \in \{0, 1\}$  and columns of **X** are linearly dependent, and **X** is called the *design matrix*
- $\triangleright$  Note that  $\mathbf{x}_j$  are known and called the explanatory variable and Y as the response variable

・ 同 ト ・ ヨ ト ・ ヨ ト

- ▷ In most cases  $x_{ij}$ s are generally chosen so that the columns of **X** are linearly independent i.e rank of **X** is d + 1
- ▷ However, in some experimental design situations  $x_{ij} \in \{0, 1\}$  and columns of **X** are linearly dependent, and **X** is called the *design matrix*
- $\triangleright$  Note that  $\mathbf{x}_j$  are known and called the explanatory variable and Y as the response variable

Question How do we estimate  $\beta$ ?

<日<br />
<</p>

- ▷ In most cases  $x_{ij}$ s are generally chosen so that the columns of **X** are linearly independent i.e rank of **X** is d + 1
- ▷ However, in some experimental design situations  $x_{ij} \in \{0, 1\}$  and columns of **X** are linearly dependent, and **X** is called the *design matrix*
- $\triangleright$  Note that  $\mathbf{x}_j$  are known and called the explanatory variable and Y as the response variable

Question How do we estimate  $\beta$ ?

Note that the first column of X need not be 1.

<日<br />
<</p>

- ▷ In most cases  $x_{ij}$ s are generally chosen so that the columns of **X** are linearly independent i.e rank of **X** is d + 1
- ▷ However, in some experimental design situations  $x_{ij} \in \{0, 1\}$  and columns of **X** are linearly dependent, and **X** is called the *design matrix*
- $\triangleright$  Note that  $\mathbf{x}_j$  are known and called the explanatory variable and Y as the response variable

Question How do we estimate  $\beta$ ?

Note that the first column of X need not be 1.

 $\begin{bmatrix} Y_1 \\ Y_2 \\ \vdots \\ Y_n \end{bmatrix} = \begin{bmatrix} x_{10} & x_{11} & x_{12} & \dots & x_{1,d-1} \\ x_{20} & x_{21} & x_{22} & \dots & x_{2,d-1} \\ \vdots & \vdots & \vdots & \vdots & \vdots \\ x_{n0} & x_{n1} & x_{n2} & \dots & x_{n,d-1} \end{bmatrix} \begin{bmatrix} \beta_0 \\ \beta_1 \\ \vdots \\ \beta_{d-1} \end{bmatrix} + \begin{bmatrix} \epsilon_1 \\ \epsilon_2 \\ \vdots \\ \epsilon_n \end{bmatrix}$