

UNIVERSITÀ DEGLI STUDI DI BERGAMO

Dipartimento di Ingegneria Gestionale, dell'Informazione e della Produzione

### Lesson 2.

### Linear regression

DATA SCIENCE AND AUTOMATION COURSE

MASTER DEGREE SMART TECHNOLOGY ENGINEERING

TEACHER Mirko Mazzoleni

PLACE University of Bergamo

### Outline

1. Components of learning

2. Linear regression

3. Least squares

#### 4. Gradient descent



### Outline

1. Components of learning

2. Linear regression

3. Least squares

#### 4. Gradient descent



A Dipartimento
 Di di Ingegneria Gestionale,
 O dell'Informazione e della Produzione

## **Components of learning (in general)**

- Input: x (e-mail textual content)  $\rightarrow$  each dimension is some e-mail attribute
- Output:  $y (\text{spam / not spam?}) \rightarrow \text{the decision that we have to take in the end}$
- Target function:  $f: \mathcal{X} \to \mathcal{Y}$  (Ideal spam filter formula)  $\to$  unknown, we have to learn it
- Data:  $\mathcal{D} = \{(x(1), y(1)), \dots, (x(N), y(N))\}$  (historical records of e-mail examples)
  - ✓ Each feature vector x consists of different regressors or features, i.e. information used to predict the output variable
- Hypothesis:  $g: \mathcal{X} \to \mathcal{Y}, g \in \mathcal{H}$  (formula to be used)  $\to g$  is an **approximation** of f

# ${\mathcal H}$ is called the ${\bf Hypothesis}\ {\bf space}.$ This, together with the ${\bf Learning}\ {\bf algorithm},$ form the ${\bf learning}\ {\bf model}$



## **Supervised learning**

- The "correct answer" (output label) y is given
- Predict y from a set of inputs  $x \in \mathbb{R}^{d \times 1}$
- **Regression:** predict a continuous output  $y \in \mathbb{R}$  (real value)
- **Classification:** predict a discrete categorical output  $y \in \{1, 2, ..., C\}$  (class)





### **Example: house prices regression** Single feature x<sub>3</sub>

Suppose we want to find a linear function which relates the measured regressors x<sub>1</sub>, x<sub>2</sub>, x<sub>3</sub>, x<sub>4</sub> with the **observed** output y
 Output

	Size $[feet^2]$	Number of bedrooms	Number of floors	Age of home [year]	Price [\$] Variable y
	2104	5	1	45	$4.60\cdot 10^5$
s N	1416	3	2	40	$2.32\cdot 10^5$
ion	1534	2	1	30	$3.15\cdot 10^5$
vat	:	:	:		
Isc	$\downarrow$	$\downarrow$	$\downarrow$	$\downarrow$	Single observation
ŏ	$x_1$	$x_2$	$x_3$	$x_4$	$\frac{y}{y}$ (feature vector) $x$

- The number of rows is the number of data points (also known as number of observations) N
- The *i*-th observation is the vector  $\mathbf{x}(i) = [x_1(i) \ x_2(i) \ x_3(i) \ x_4(i)]^T \in \mathbb{R}^{4x_1}$
- Each feature vector x has associated a response  $y \in \mathbb{R}$  that we want to predict for new observations  $x^*$



Number of

### Example: house prices classification

- The components of the features vector are the same. The difference lies in the response variable, which now is a **class** (categorical data type) and not a real value
- Suppose that instead of the price value in dollars, we want to classify houses as **expensive** (class y = 1) or **cheap** (class y = 0)

Size $[feet^2]$	Number of bedrooms	Number of floors	Age of home $\left[\mathrm{year}\right]$	Price [class]
2104	5	1	45	1
1416	3	2	40	0
1534	2	1	30	1
÷	:	:	:	÷
$\downarrow$	$\downarrow$	$\downarrow$	$\downarrow$	$\downarrow$
$x_1$	$x_2$	$x_3$	$x_4$	$  \qquad y$

• The point x could be classified to class y = 1 if the probability of x to belong to class 1 is  $\ge 0.5$ 



### **Unsupervised learning**

- Instead of (input, output) we get (input, ?)
- Here there is no a function *f* to learn
- Find properties of the inputs  $x \in \mathbb{R}^{d \times 1}$
- High-level representation of the input
- Elements into the same cluster have similar properties



UNIT DEG DI B

### **Reinforcement learning**

- Instead of (input, output) we get (input, output, reward)
- The algorithm tries to learn what action to take, in order to maximize the reward
- This is called a policy
- Applications in control, robotics, A/B testing





### Business problems as data science examples - revisited

#### **Supervised**

#### Unsupervised

- Spam e-mail detection system
   Classification
- Credit approval Classification
- Recognize objects in images **Classification**
- Find the relation between house prices and house sizes **Regression**
- Predict the stock market Regression

- Market segmentation Clustering
- Language models (word2vec)
   Similarity matching

### Social network analysis Link

**Data reduction** 

- Low-order data representations
- Movies recommendation

Similaritiy matching

#### **Supervised or unsupervised**



### **Supervised learning: problem statement**

The aim is to **learn an unknown function** f given a dataset  $\mathcal{D}$ 

- The function is searched in the hypothesis space  $\mathcal{H}$ , where  $h \in \mathcal{H}$  is a specific function
- We want to find a function h that approximates f well, on the **whole domain**  $\mathcal{X}$

What does  $h \approx f$  mean?

- We need to define an error measure
- Almost always this is a **pointwise** definition: e(f(x), h(x))



### **Cost functions**

#### **Pointwise error examples**

- Squared error:  $e(f(x), h(x)) = (f(x) h(x))^2 \rightarrow \text{used for regression}$
- Binary error:  $e(f(x), h(x)) = \mathbb{I}[f(x) \neq h(x)] \rightarrow \text{used for classification}$

It is interesting to look at the **overall error**, which considers all *N* examples:

#### **Overall error examples**

- In-sample error:  $E_{in} = \frac{1}{N} \sum_{i=1}^{N} e(f(\mathbf{x}), h(\mathbf{x})) \rightarrow$  error on data that I actually have
- Out-of-sample error:  $E_{out} = \mathbb{E}_x[e(f(x), h(x))] \rightarrow \text{error on data I could possibly observe}$



### Outline

1. Components of learning

#### 2. Linear regression

3. Least squares

#### 4. Gradient descent



À | Dipartimento
 Di | di Ingegneria Gestionale,
 O | dell'Informazione e della Produzione

### Linear regression

**Aim:** Suppose to have at disposal a dataset  $\mathcal{D} = \{(x(1), y(1)), ..., (x(N), y(N))\}$ . Find the relation between a set of input variables  $x \in \mathbb{R}^{(d-1)\times 1}$  and an output variable  $y \in \mathbb{R}$ , using

- The vector  $\theta$  is called **parameters vector**  $\rightarrow$  to be found by minimizing a cost function
- The vector x(i) is called **features vector** for *i*-th observation  $\rightarrow$  attributes of individuals
- The quantity  $\epsilon(i)$  is the error due not perfect explanation of the y(i) using x(i)



### **Geometrical interpretation**

#### **One-variable case**

In this case, there is only **one feature**  $x_1$ and **two parameters**  $\theta_0, \theta_1$ 

 $y(i) = \theta_0 + \theta_1 x_1(i) + \epsilon(i)$ 



#### Two-variables case

In this case, there are **two features**  $x_{1,}x_{2}$ and **three parameters**  $\theta_{0}, \theta_{1}, \theta_{2}$ 

$$y(i) = \theta_0 + \theta_1 x_1(i) + \theta_2 x_2 + \epsilon(i)$$





### Outline

1. Components of learning

2. Linear regression

#### 3. Least squares

#### 4. Gradient descent



À Dipartimento
 DI di Ingegneria Gestionale,
 dell'Informazione e della Produzione

### Least squares cost function

How well does  $h(x) = x^T \theta$  approximates f(x)?

 Linear regression estimates the parameters θ using the least squares method, i.e. by minimizing the squared error between observed and predicted output





### **Cost function - minimization**

$$J(\boldsymbol{\theta}) = E_{in}(\boldsymbol{\theta}) = \frac{1}{N} \sum_{i=1}^{N} (\boldsymbol{y(i)} - \boldsymbol{x(i)}^T \boldsymbol{\theta})^2$$

- Since the model is linear in the parameters and the error measure is quadratic, the cost function is convex → it admits a unique (global) minimum
- In this case the minimum can be found in **closed-form**



### Least squares cost function – matrix form

We can express the linear regression problem using matrices Features vector  $x^{T}(i)$  $X = \begin{bmatrix} 1 & x_1(1) & x_2(1) & \cdots & x_{d-1}(1) \\ 1 & x_1(2) & x_2(2) & & x_{d-1}(2) \\ \vdots & \vdots & & \ddots & \vdots \\ 1 & x_1(N) & x_2(N) & \cdots & x_{d-1}(N) \end{bmatrix} \quad \boldsymbol{\theta} = \begin{bmatrix} \theta_0 \\ \theta_1 \\ \vdots \\ \theta_{d-1} \end{bmatrix} \quad \begin{array}{c} Y = \begin{bmatrix} y(1) \\ y(2) \\ \vdots \\ N \times 1 \end{bmatrix} \quad \begin{array}{c} E \\ E \\ y(N) \end{bmatrix} \quad \begin{array}{c} \varepsilon(1) \\ \varepsilon(2) \\ \vdots \\ \varepsilon(N) \end{bmatrix}$  $\Gamma = T (1) T$ 

$$X_{N \times d} = \begin{bmatrix} \mathbf{x}^{T} (1) \\ \mathbf{x}^{T} (2) \\ \vdots \\ \mathbf{x}^{T} (N) \end{bmatrix} \qquad Y = X\mathbf{\theta} + E \Rightarrow \begin{bmatrix} J(\mathbf{\theta}) = \frac{1}{N} \|Y - X\mathbf{\theta}\|_{2}^{2} = \frac{1}{N} (Y - X\mathbf{\theta})^{T} (Y - X\mathbf{\theta}) \\ 1 \times N \end{bmatrix}$$



A | Dipartimento
 di Ingegneria Gestionale,
 dell'Informazione e della Produzione

### Least squares cost function – matrix form

It is useful to remember these matrix derivation properties:

 $\nabla_{\mathbf{x}}(\mathbf{x}^T \cdot A \cdot \mathbf{x}) = (A + A^T) \cdot \mathbf{x}$ 1×d d×d d×1 d×d d×1  $\nabla_{\mathbf{x}}(\mathbf{x}^T \cdot \mathbf{b}) = \mathbf{b}$ 1×d d×1 d×1 d×1 d×1  $1 \times d \quad d \times 1 \quad d \times 1$  $J(\boldsymbol{\theta}) = \frac{1}{N} (Y - X\boldsymbol{\theta})^T (Y - X\boldsymbol{\theta}) = \frac{1}{N} (Y^T Y - Y^T X \cdot \boldsymbol{\theta}_{1 \times N N \times d d \times 1} - \boldsymbol{\theta}^T \cdot X^T Y + \boldsymbol{\theta}^T \cdot X^T X \cdot \boldsymbol{\theta})$  $= \frac{1}{M} (Y^T Y - 2 \cdot \boldsymbol{\theta}^T \cdot X^T Y + \boldsymbol{\theta}^T \cdot X^T X \cdot \boldsymbol{\theta})$  $\nabla J(\boldsymbol{\theta}) = \mathbf{0} \Rightarrow \frac{1}{N} (-2X^T Y + 2X^T X \boldsymbol{\theta}) = \mathbf{0} \Rightarrow \qquad \widehat{\boldsymbol{\theta}} = (X^T X)^{-1} X^T Y$ 



### **Normal equations**

$$\widehat{\boldsymbol{\theta}} = (X^T X)^{-1} X^T Y$$
 Normal equations

What if the matrix  $X^T X$  is **not invertible**?  $\longrightarrow$  Use **pseudo-inverse**. In MatLab:

- Redundant features (linearly dependent)
  - $\checkmark x_1 = \text{height in m}^2$
  - ✓  $x_2$  = height in feet<sup>2</sup>
- Too many features (ex.  $N \leq d$ )
  - ✓ Delete some feature
  - $\checkmark$  Use regularization (later in the course)

theta hat = pinv(X' \* X) \* X \* y

- The method of normal equations is
   **slow** if d is very large
  - $\checkmark~$  To solve this, iterative methods
    - using gradient descent can be

used



### Outline

1. Components of learning

2. Linear regression

3. Least squares

#### 4. Gradient descent



### **Gradient descent**

• The **gradient descent** is a general iterative method for minimizing differentiable functions

• The value of the parameters at iteration t + 1 is (given a random initial point  $\hat{\theta}(0)$ )

$$\widehat{\boldsymbol{\theta}}(t+1) = \widehat{\boldsymbol{\theta}}(t) - \alpha \cdot \nabla J(\boldsymbol{\theta}) \Big|_{\boldsymbol{\theta} = \widehat{\boldsymbol{\theta}}(t)}$$

 $\alpha \in \mathbb{R}_{>0}$ : learning rate



di Ingegneria Gestionale, di Ingegneria Gestionale, dell'Informazione e della Produzione

### **Gradient descent**

Consider the **scalar** case where  $\theta \in \mathbb{R}$ 

$$\hat{\theta}(t+1) = \hat{\theta}(t) - \alpha \nabla J(\theta) \Big|_{\theta = \hat{\theta}(t)}$$

• 
$$\nabla J(\theta)|_{\theta=\widehat{\theta}(t)} > 0 \Rightarrow \widehat{\theta}(t+1) < \widehat{\theta}(t)$$

# The new estimate is closer to the optimal value $\theta^*$





### **Gradient descent**

Consider the **scalar** case where  $\theta \in \mathbb{R}$ 

$$\hat{\theta}(t+1) = \hat{\theta}(t) - \alpha \nabla J(\theta) \Big|_{\theta = \hat{\theta}(t)}$$

• 
$$\nabla J(\theta)|_{\theta=\widehat{\theta}(t)} < 0 \Rightarrow \widehat{\theta}(t+1) > \widehat{\theta}(t)$$

# The new estimate is closer to the optimal value $\theta^*$





### Linear regression laboratory: predict house prices

Т

We want to **predict the price** of the houses in Portland, Oregon

- Each house is described by the following features
   ✓ x<sub>1</sub>: Number of **bedrooms** ✓ x<sub>2</sub>: Size [feet<sup>2</sup>]
- The **training set** consists of N = 47 houses with  $x_1(i), x_2(i)$  and y(i), for i = 1, ..., N

$$y(i) = \mathbf{x}(i)^{T} \boldsymbol{\theta} + \epsilon(i) \qquad \mathbf{x}(i) = \begin{bmatrix} 1 & x_{1}(i) & x_{2}(i) \end{bmatrix}$$
$$X = \begin{bmatrix} \mathbf{x}^{T}(1) \\ \mathbf{x}^{T}(2) \\ \vdots \\ \mathbf{x}^{T}(N) \end{bmatrix} \qquad \boldsymbol{\theta} = \begin{bmatrix} \theta_{0} \\ \theta_{1} \\ \theta_{2} \end{bmatrix} \qquad Y = \begin{bmatrix} y(1) \\ y(2) \\ \vdots \\ 47 \times 1 \end{bmatrix} \begin{bmatrix} y(1) \\ y(2) \\ \vdots \\ y(47) \end{bmatrix}$$

```
% Read data from file
data = csvread(`housedata.txt');
X = data(:, 1:2); % Features
y = data(:, 3); % Price
N = length(y); % Number of data
```

```
% Add intercept term to X
X = [ones(N, 1) X];
```

```
% Calculate the parameters from
the normal equation
theta_hat = pinv(X'*X)*X'*y;
% Estimate the price of a 1650
sq-ft, 3 br house
price_hat = [1 3 1650]*theta_hat;
Point not seen during training (i.e. estimation of 0)
```

