

# Linear & Logistic Regression

FTI UII

13 Januari 2020

Dr. Ing. Ridho Rahmadi, S.Kom., M.Sc.

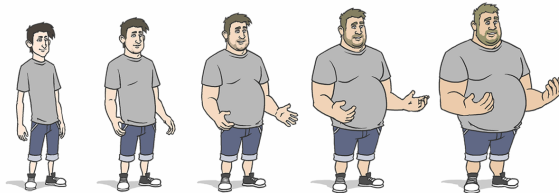


UNIVERSITAS  
ISLAM  
INDONESIA

- Head of Center of Data Science, UII
- Head of Research Laboratory, Informatics, UII
- Education
  - Universitas Islam Indonesia (S1)
  - Czech Technical University in Prague (S2)
  - Johannes Kepler University, Austria (S2)
  - Radboud University Nijmegen, the Netherlands (S3)
  - Carnegie Mellon University, USA (Visiting scholar)
- Research interest
  - Machine learning, deep learning, causal modeling, stability selection, multi-objective evolutionary algorithms
- [ridho.rahmadi@uii.ac.id](mailto:ridho.rahmadi@uii.ac.id)
- 081129513045

# How do we see the world?

2 | 53

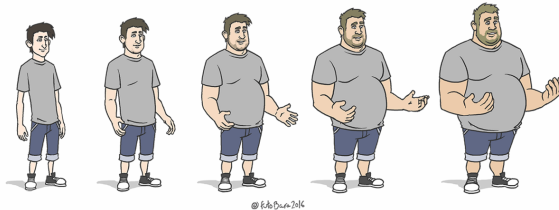


© K. A. D. 2016



# How do we see the world?

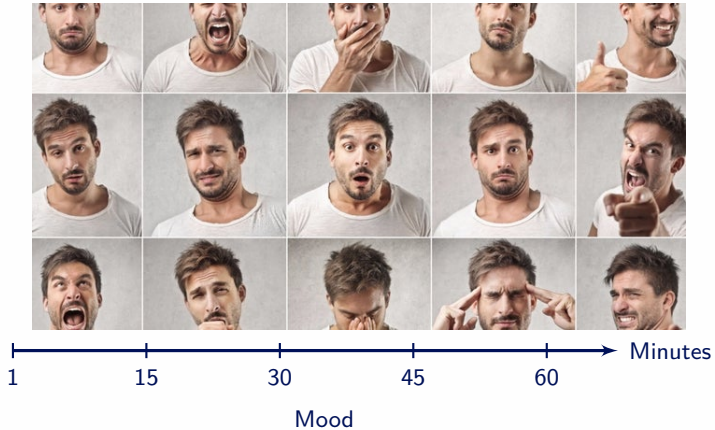
2 | 53



The more cookies I eat, the more weight I gain

# How do we see the world?

3 | 53



# How do we see the world?

3 | 53



I am not necessarily always happier as time passes

## How do we *see* the world?

4 | 53

- If I eat 1 cookie per day, I gain 2 Kg weight
- If I eat 2 cookies per day, I gain 4 Kg weight

## How do we *see* the world?

4 | 53

- If I eat 1 cookie per day, I gain 2 Kg weight
- If I eat 2 cookies per day, I gain 4 Kg weight

The more cookies I eat, the more weight I gain.



## How do we see the world?

4 | 53

- If I eat 1 cookie per day, I gain 2 Kg weight
- If I eat 2 cookies per day, I gain 4 Kg weight

The more cookies I eat, the more weight I gain.

We call such an example as a **linear model**.

## How do we see the world?

4 | 53

- If I eat 1 cookie per day, I gain 2 Kg weight
- If I eat 2 cookies per day, I gain 4 Kg weight

The more cookies I eat, the more weight I gain.

We call such an example as a **linear model**. For example,

- If A goes up, so does B, OR
- If A goes up, B goes down

## How do we see the world?

4 | 53

- If I eat 1 cookie per day, I gain 2 Kg weight
- If I eat 2 cookies per day, I gain 4 Kg weight

The more cookies I eat, the more weight I gain.

We call such an example as a **linear model**. For example,

- If A goes up, so does B, OR
- If A goes up, B goes down

Given two variables, a linear relationship among them indicates consistent directions of changes.

## How do we *see* the world?

5 | 53

I just woke up and

- I am happy (minute 1)
- I am angry as no one WA me (minute 10)
- I am happy as 1 WA comes (minute 11)
- I am nervous as I'll have an exam (minute 30)

## How do we see the world?

5 | 53

I just woke up and

- I am happy (minute 1)
- I am angry as no one WA me (minute 10)
- I am happy as 1 WA comes (minute 11)
- I am nervous as I'll have an exam (minute 30)

I am not necessarily always happier as time passes.

## How do we see the world?

5 | 53

I just woke up and

- I am happy (minute 1)
- I am angry as no one WA me (minute 10)
- I am happy as 1 WA comes (minute 11)
- I am nervous as I'll have an exam (minute 30)

I am not necessarily always happier as time passes.

We call such an example as **nonlinear model**.

## How do we see the world?

5 | 53

I just woke up and

- I am happy (minute 1)
- I am angry as no one WA me (minute 10)
- I am happy as 1 WA comes (minute 11)
- I am nervous as I'll have an exam (minute 30)

I am not necessarily always happier as time passes.

We call such an example as **nonlinear model**.

For example, if A goes up, B alternates up and down.

## How do we see the world?

6 | 53

- If I eat 1 cookie per day, I gain 2 Kg of weight
- If I eat 2 cookies per day, I gain 4 Kg of weight

What if I eat 3 cookies? How many Kg of weight will I gain?





## How do we see the world?

7 | 53

I just woke up and

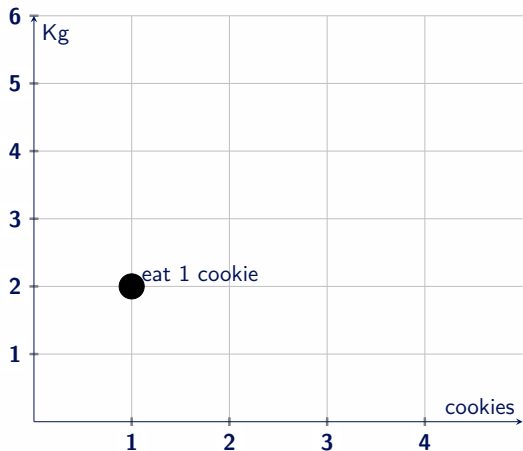
- I am happy (minute 1)
- I am angry as no one WA me (minute 10)
- I am happy as 1 WA comes (minute 11)
- I am nervous as I'll have an exam (minute 30)

Will I be happy in minute 110? Or nervous?



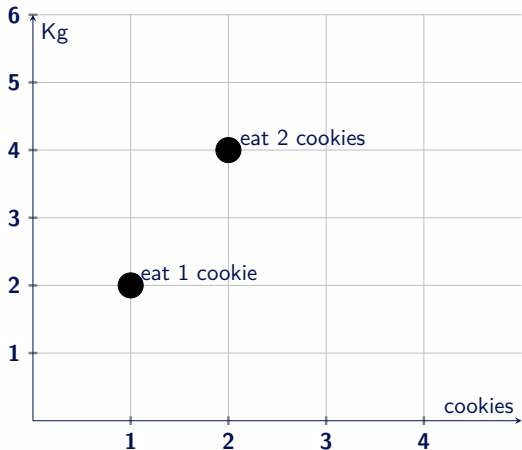
## Say it in Mathematics

8 | 53



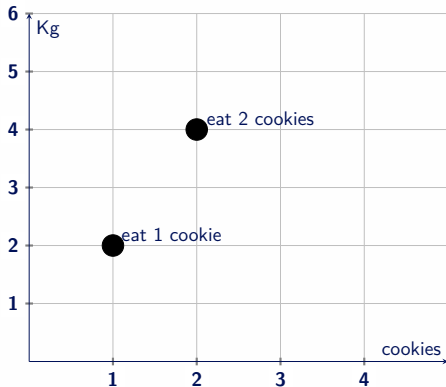
## Say it in Mathematics

9 | 53



## Say it in Mathematics

10 | 53



What if I eat 3 cookies?

## Say it in Mathematics

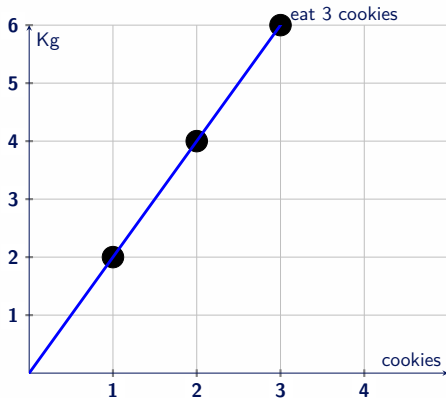
11 | 53



We can draw a **line** that passes through the two points.

## Say it in Mathematics

12 | 53



Using the line, we can predict “what if I eat 3 cookies?”

## Say it in Mathematics

13 | 53



With the line, you can predict the weight gain for any (positive) number of cookies eaten.

In terms of mathematical function, a line can be represented by

$$y = \theta_0 + \theta_1 x,$$

where  $\theta_0$  is the **intercept** and  $\theta_1$  is the **slope**.



In terms of mathematical function, a line can be represented by

$$y = \theta_0 + \theta_1 x,$$

where  $\theta_0$  is the **intercept** and  $\theta_1$  is the **slope**.

The intercept indicates the **y**-coordinate through which the line passes. It is the value of **y** when **x** = **0**.

In terms of mathematical function, a line can be represented by

$$y = \theta_0 + \theta_1 x,$$

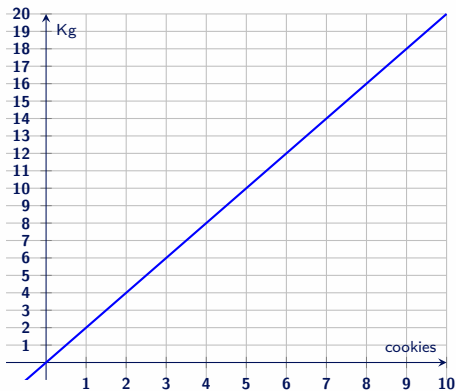
where  $\theta_0$  is the **intercept** and  $\theta_1$  is the **slope**.

The intercept indicates the **y**-coordinate through which the line passes. It is the value of **y** when **x = 0**.

The slope represents how steep the line is.

## Say it in Mathematics

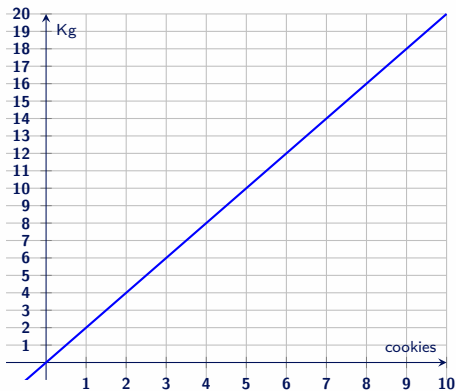
15 | 53



From the previous example, we have  $y = 0 + 2x$ .

## Say it in Mathematics

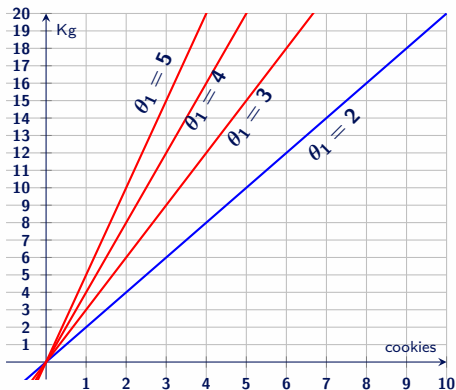
15 | 53



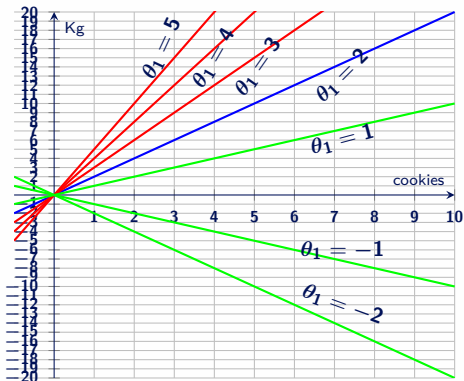
From the previous example, we have  $y = 0 + 2x$ . Now let's play a bit with the slope and intercept.

## Say it in Mathematics

16 | 53



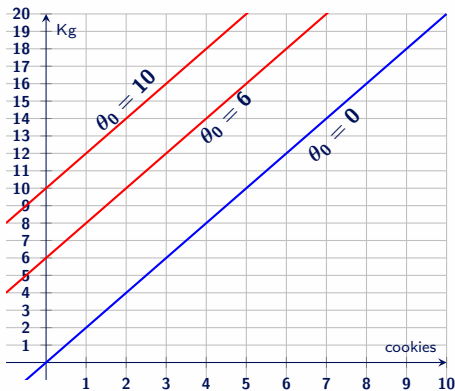
From  $y = 0 + 2x$ , we have  $\theta_1 = 2$ . See what happens if we **increase** the slope  $\theta_1$ .



From  $y = 0 + 2x$ , we have  $\theta_1 = 2$ . See what happens if we **decrease** the slope  $\theta_1$ .

## Say it in Mathematics

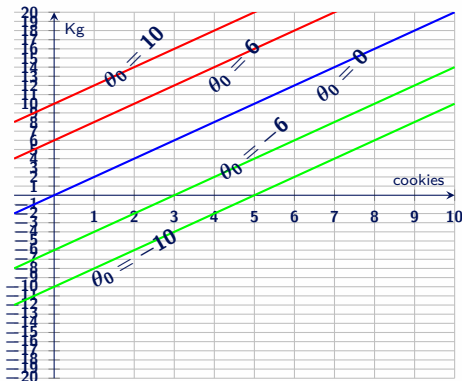
18 | 53



From  $y = 0 + 2x$ , we have  $\theta_0 = 0$ . See what happens if we **increase** the intercept  $\theta_0$ .

## Say it in Mathematics

19 | 53



From  $y = 0 + 2x$ , we have  $\theta_0 = 0$ . See what happens if we **decrease** the intercept  $\theta_0$ .



The previous example illustratively describes the idea behind the **linear regression**.

That is, based on the data we have, we want to predict the weight gained, given the number of cookies we eat, by fitting a line.

Next, we will describe the linear regression in more detail.

## Real-world cases

21 | 53

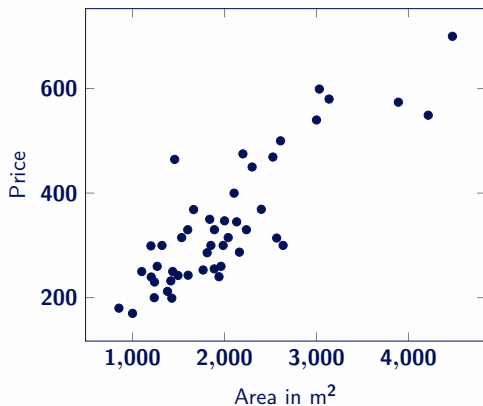
In real-world cases, the data sets are often of the form  $(\mathbf{x}^{(i)}, \mathbf{y}^{(i)}); i = 1, \dots, m$ . For example,

House area ( $\mathbf{x}$ )	Price ( $\mathbf{y}$ )
2104	400
1600	330
2400	369
1416	232
3000	540
$\vdots$	$\vdots$

Based on the data above, a typical question is, e.g., what is the price of a house if the area is 558 M<sup>2</sup>?

## Plot the data

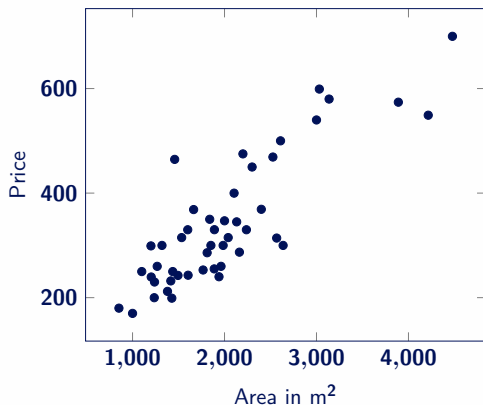
22 | 53



Does the data distribution form a linear fashion?

## Plot the data

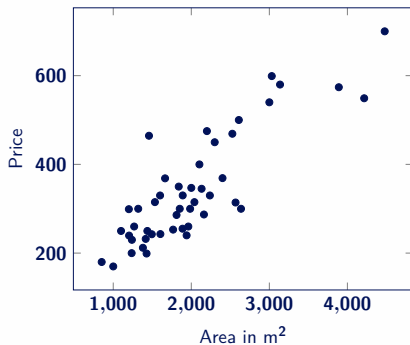
23 | 53



A typical question is, e.g., what is the price of a house if the area is 558 M²?

## Plot the data

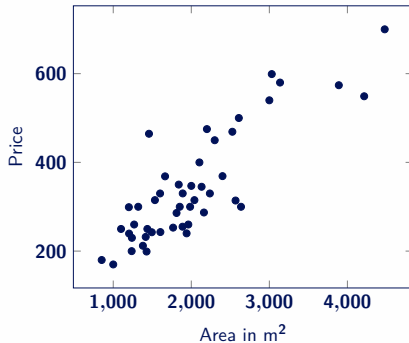
24 | 53



- Note that we are interested in to predict *unseen*  $\hat{y}$  based on  $\hat{x}$  from other population

## Plot the data

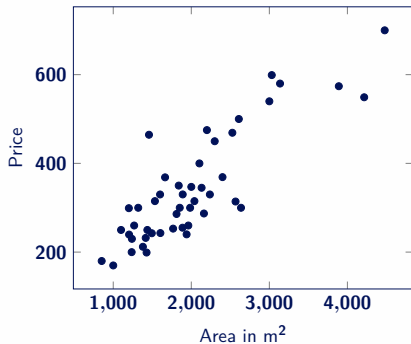
24 | 53



- Note that we are interested in to predict *unseen*  $\hat{y}$  based on  $\hat{x}$  from other population
- Solution: we can find a line (model) that constitutes the data we have, and use the line to predict  $\hat{y}$  based on  $\hat{x}$

## Plot the data

24 | 53

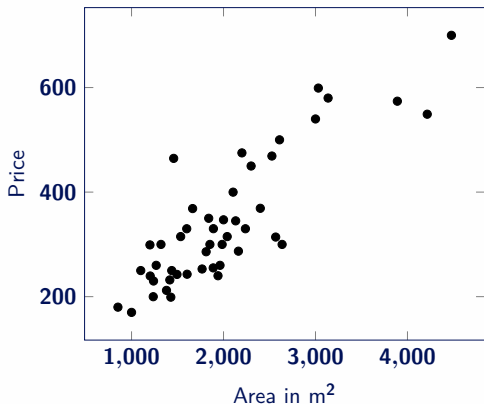


- Note that we are interested in to predict *unseen*  $\hat{y}$  based on  $\hat{x}$  from other population
- Solution: we can find a line (model) that constitutes the data we have, and use the line to predict  $\hat{y}$  based on  $\hat{x}$

This is called a model **generalization**, i.e., you obtain a model from a data set and apply it to another data set(s). This is a fundamental concept in Machine & Deep Learning.

## Linear regression

25 | 53

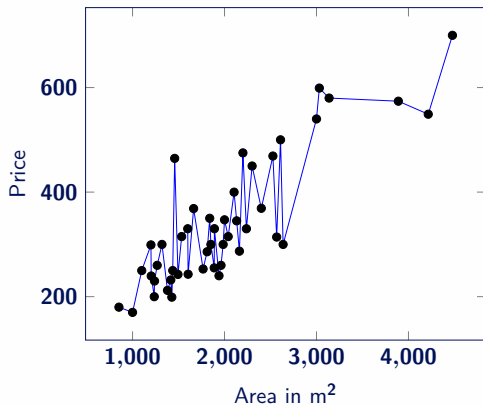


The data do form a linear fashion, but how to find a good line/model?



## Linear regression

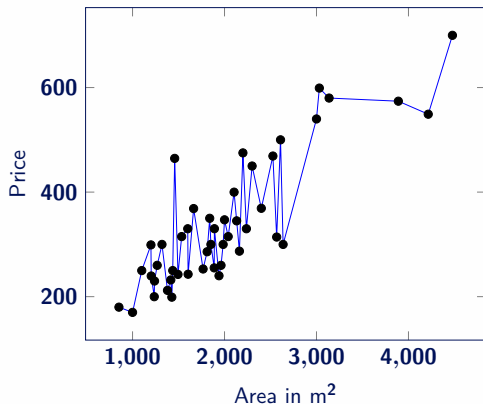
26 | 53



Draw a line by connecting all the points like this?

## Linear regression

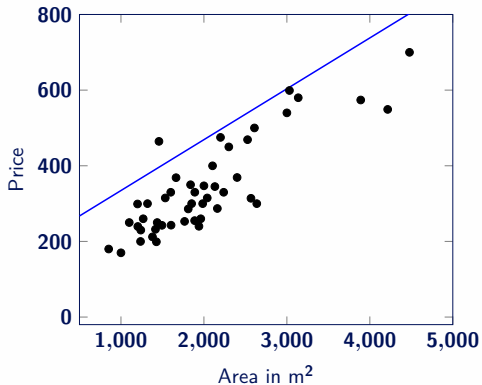
26 | 53



Draw a line by connecting all the points like this? Recall that our objective is a model generalization; the model above will not fit other data.

## Linear regression

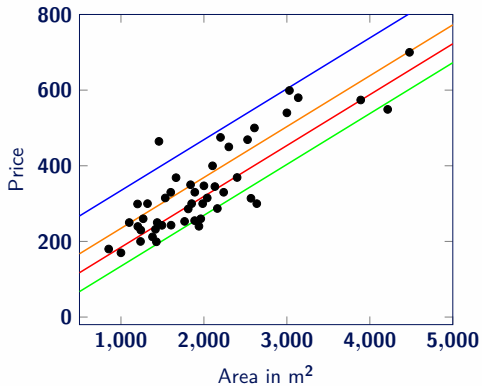
27 | 53



Draw a line like this?

## Linear regression

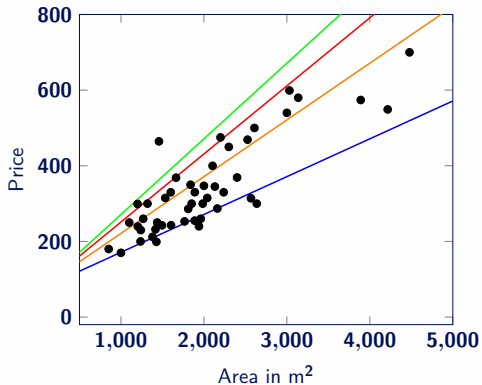
28 | 53



But which line?

## Linear regression

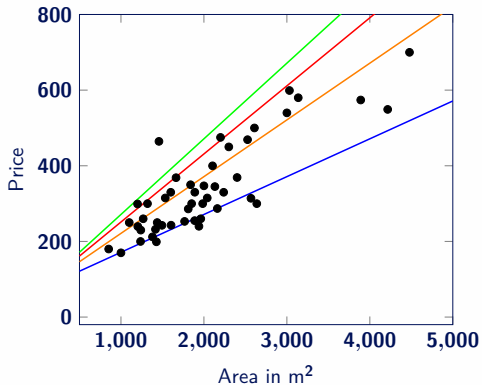
29 | 53



But which line?

## Linear regression

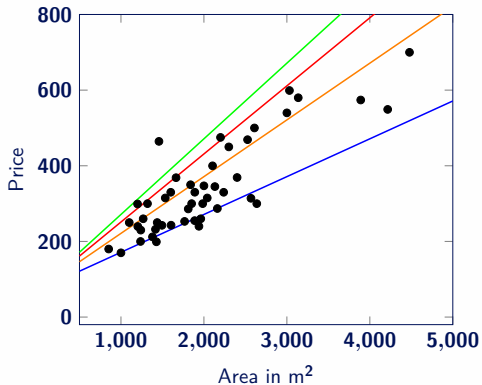
29 | 53



But which line? What is the criteria of a *good* line/model? Define ones.

# Linear regression

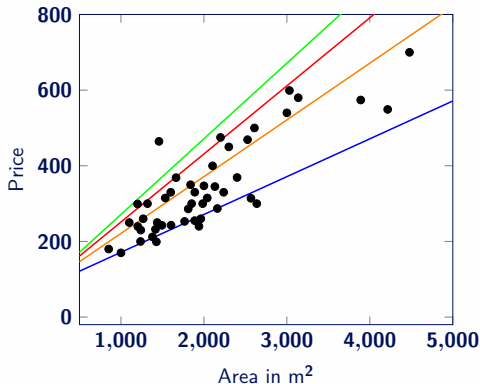
30 | 53



Informally, we can think of a *good* line/model is the one that is generally close to all data points.

# Linear regression

31 | 53

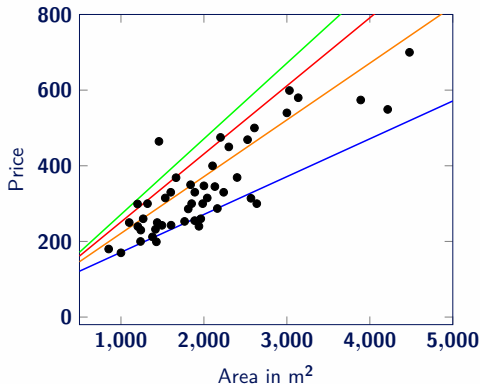


To indicate how close a line, we can measure distances between data points to the line. The total distance is often called **error**; the lower it is, the better.



# Linear regression

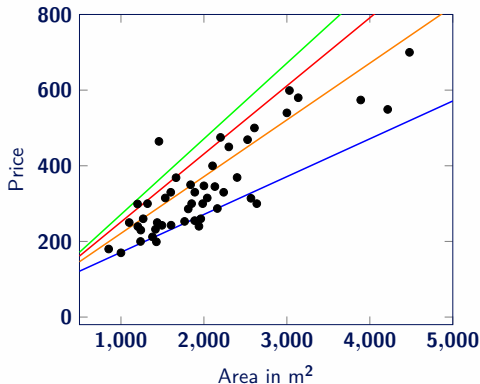
32 | 53



Linear regression is about to find the *best* line by selecting the one with the minimum error. How?

## Linear regression

33 | 53



Recall that we can search lines by changing the values of intercept  $\theta_0$  and slope  $\theta_1$  in  $y = \theta_0 + \theta_1 x$ . But of course, we do not want to search **randomly**.

## Linear Regression

34 | 53

More formally, a straight line can be represented by,

$$h(\mathbf{x}) = \sum_{i=0}^n \theta_i x_i = \boldsymbol{\theta}^T \mathbf{x},$$

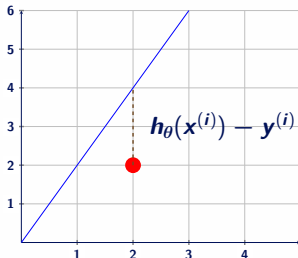
# Linear Regression

34 | 53

More formally, a straight line can be represented by,

$$h(\mathbf{x}) = \sum_{i=0}^n \theta_i x_i = \boldsymbol{\theta}^T \mathbf{x},$$

- The error term indicates the distance between a point  $\mathbf{y}^{(i)}$  to the line  $\mathbf{h}_{\boldsymbol{\theta}}(\mathbf{x}^{(i)})$ , i.e.,  $\mathbf{h}_{\boldsymbol{\theta}}(\mathbf{x}^{(i)}) - \mathbf{y}^{(i)}$ .



The total error, or often called **cost function** in ML/DL, thus can be written

$$J(\theta) = \frac{1}{2} \sum_{i=1}^m (h_{\theta}(\mathbf{x}^{(i)}) - \mathbf{y}^{(i)})^2.$$

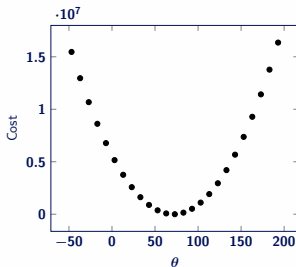
# Linear Regression

35 | 53

The total error, or often called **cost function** in ML/DL, thus can be written

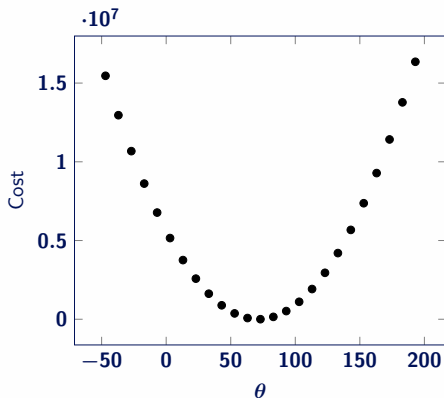
$$J(\theta) = \frac{1}{2} \sum_{i=1}^m (h_{\theta}(x^{(i)}) - y^{(i)})^2.$$

The function above is called the **ordinary least square**. Here is the plot.



## Ordinary least square

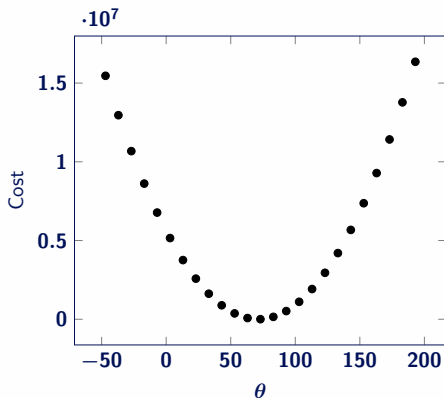
36 | 53



- Hey, we get some clue here!

## Ordinary least square

36 | 53

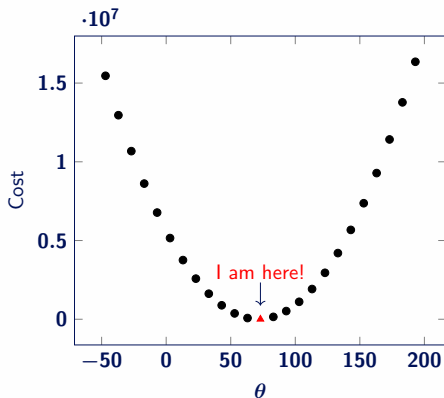


- Hey, we get some clue here!
- Obviously we know which one the minimum cost is, don't we?



# Ordinary least square

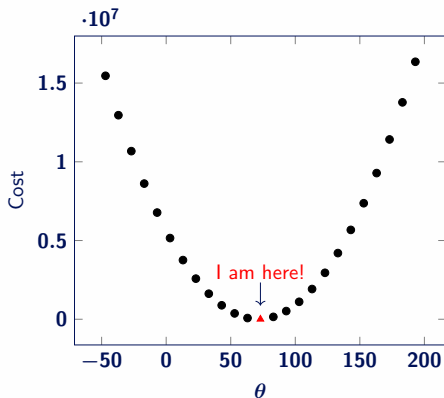
37 | 53



- Hey, we get some clue here!
- Obviously we know which one the minimum cost is, don't we?
- That is, look at "I am here"

# Ordinary least square

37 | 53

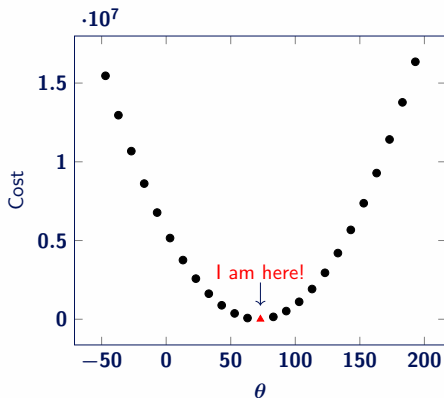


- Hey, we get some clue here!
- Obviously we know which one the minimum cost is, don't we?
- That is, look at "I am here"

Each cost (data point) above represents an error term of a line/model. The question: **How to find the minimum error?**

## Ordinary least square

38 | 53

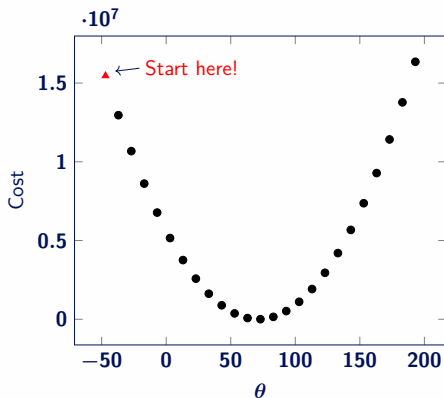


- Hey, we get some clue here!
- Obviously we know which one the minimum cost is, don't we?
- That is, look at "I am here"

How to find the minimum error? Answer: Simply going down to the *valley bottom*.

## Ordinary least square

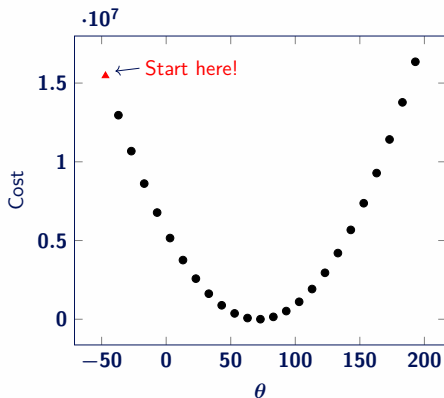
39 | 53



- Typically your randomly initialized line returns the “start here” cost. How can we go down?

# Ordinary least square

39 | 53

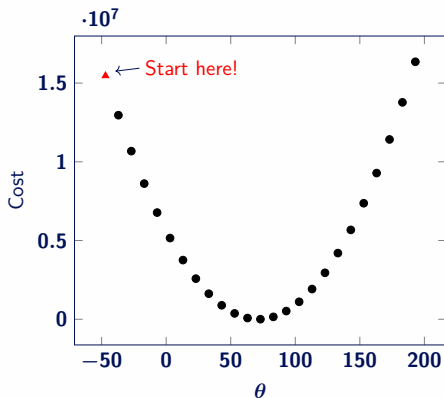


- Typically your randomly initialized line returns the “start here” cost. How can we go down?
- We can use the **gradient descent** to update **parameter  $\theta$** , so as to get the minimum cost

We now use the term **model parameter** to represent slope and intercept.

# Gradient descent

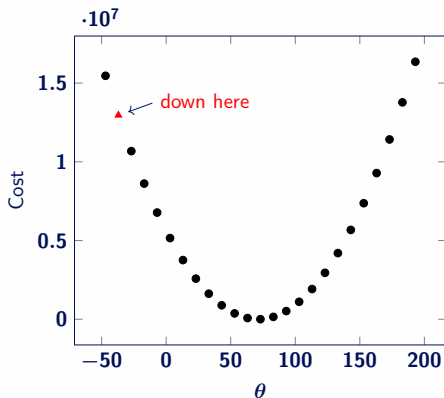
40 | 53



- Initially *guess*  $\theta$ , compute the cost (see **Start here**)

# Gradient descent

41 | 53



- Initially *guess*  $\theta$ , compute the cost (see **Start here**)
- Repeatedly, update the parameter (see **“down here”**) via

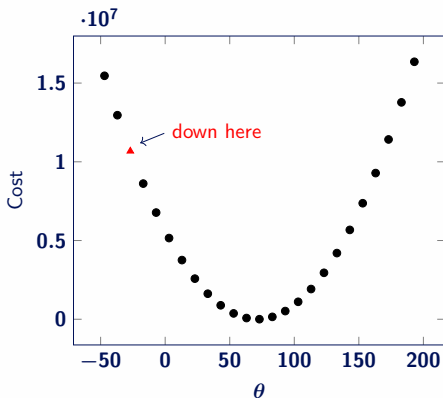
$$\theta_j := \theta_j - \alpha \frac{\partial}{\partial \theta_j} J(\theta)$$

simultaneously for all  
 $j = 0, \dots, n$ .

down again...

# Gradient descent

42 | 53



- Initially *guess*  $\theta$ , compute the cost (see **Start here**)
- Repeatedly, update the parameter (see **"down here"**) via

$$\theta_j := \theta_j - \alpha \frac{\partial}{\partial \theta_j} J(\theta)$$

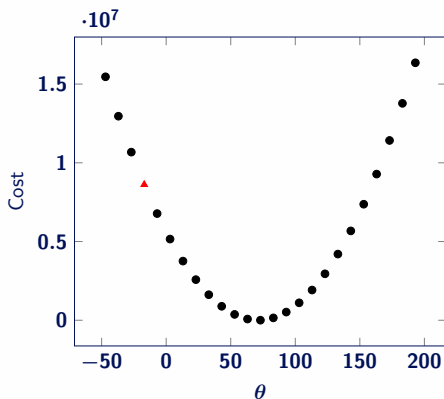
simultaneously for all  $j = 0, \dots, n$ .

down again...



# Gradient descent

43 | 53



- Initially *guess*  $\theta$ , compute the cost (see **Start here**)
- Repeatedly, update the parameter (see **“down here”**) via

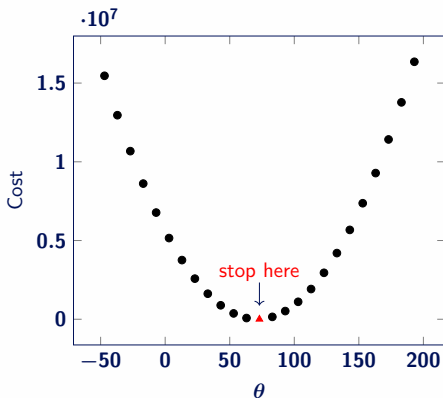
$$\theta_j := \theta_j - \alpha \frac{\partial}{\partial \theta_j} J(\theta)$$

simultaneously for all  
 $j = 0, \dots, n$ .

so on...

# Gradient descent

44 | 53



- Initially *guess*  $\theta$ , compute the cost (see **Start here**)
- Repeatedly, update the parameter (see **“down here”**) via

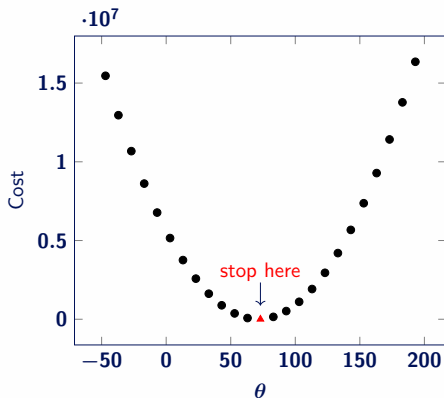
$$\theta_j := \theta_j - \alpha \frac{\partial}{\partial \theta_j} J(\theta)$$

simultaneously for all  
 $j = 0, \dots, n$ .

until reaching the minimum error.

# Gradient descent

45 | 53



- Once you reach the minimum error, the corresponding parameter  $\theta$  should give you the *fittest* line/model

1. Pick an initial line/model  $\mathbf{h}(\boldsymbol{\theta})$  by randomly choosing parameter  $\boldsymbol{\theta}$
2. Compute the corresponding cost function, e.g.,

$$J(\boldsymbol{\theta}) = \frac{1}{2} \sum_{i=1}^m (\mathbf{h}_{\boldsymbol{\theta}}(\mathbf{x}^{(i)}) - \mathbf{y}^{(i)})^2,$$

3. Update the line/model  $\mathbf{h}(\boldsymbol{\theta})$  by updating  $\boldsymbol{\theta}$  that makes  $J(\boldsymbol{\theta})$  smaller, using, e.g., the gradient descent that reads

$$\theta_j := \theta_j - \alpha \frac{\partial}{\partial \theta_j} J(\boldsymbol{\theta}),$$

where  $\alpha$  is the learning rate.

4. Repeat steps 2 and 3 until converges

## Linear regression: a summary

47 | 53

In more detail

The learning procedure above is a typical approach in mostly (parametric) models of machine and deep learning.

Note that we have been so far assuming the  $\mathbf{y}$  in our data set  $(\mathbf{x}^{(i)}, \mathbf{y}^{(i)})$ ;  $i = 1, \dots, m$  is continuous, e.g., house price, height, temperature, etc.

## Linear regression: a summary

47 | 53

In more detail

The learning procedure above is a typical approach in mostly (parametric) models of machine and deep learning.

Note that we have been so far assuming the  $\mathbf{y}$  in our data set  $(\mathbf{x}^{(i)}, \mathbf{y}^{(i)}); i = 1, \dots, m$  is continuous, e.g., house price, height, temperature, etc.

What if  $\mathbf{y}$  is discrete? E.g., pass/fail, yes/true, 1/0, etc, that leads to nonlinear fashion and a classification task.

## Linear regression: a summary

47 | 53

In more detail

The learning procedure above is a typical approach in mostly (parametric) models of machine and deep learning.

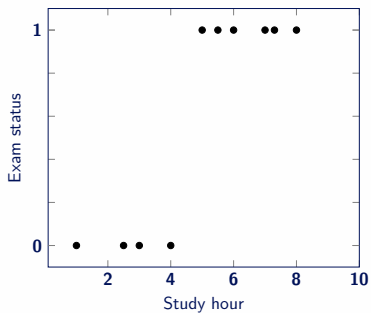
Note that we have been so far assuming the  $\mathbf{y}$  in our data set  $(\mathbf{x}^{(i)}, \mathbf{y}^{(i)}); i = 1, \dots, m$  is continuous, e.g., house price, height, temperature, etc.

What if  $\mathbf{y}$  is discrete? E.g., pass/fail, yes/true, 1/0, etc, that leads to nonlinear fashion and a classification task.

Let's see the plot example.

If  $y$  is a discrete variable

48 | 53

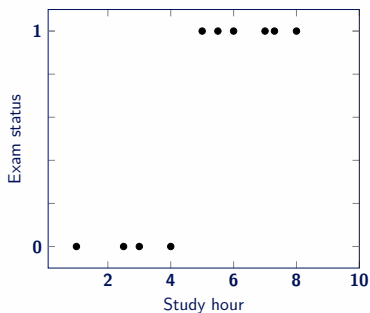


Exam status: 0 (fail), 1 (pass)



If  $y$  is a discrete variable

48 | 53

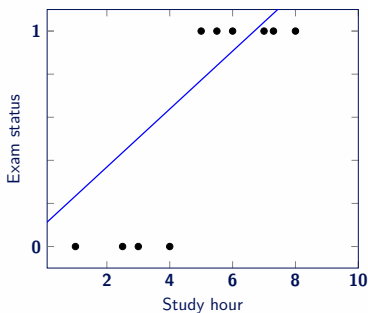


- Do you think the linear model  $y = \theta_0 + \theta_1 * x$  will fit well?

Exam status: 0 (fail), 1 (pass)

If  $y$  is a discrete variable

49 | 53

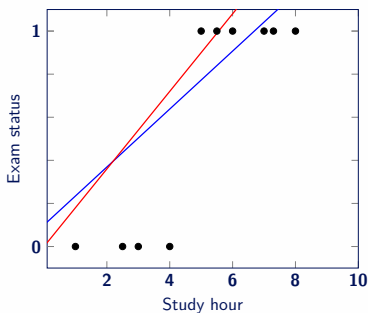


- Do you think the linear model  $y = \theta_0 + \theta_1 * x$  will fit well?
- Let's try one.

Exam status: 0 (fail), 1 (pass)

## If $y$ is a discrete variable

50 | 53

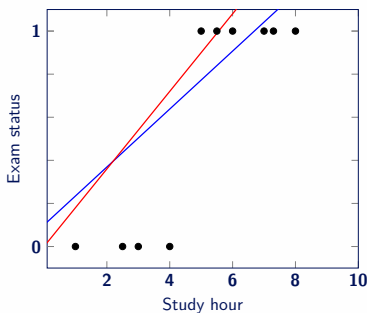


Exam status: 0 (fail), 1 (pass)

- Do you think the linear model  $y = \theta_0 + \theta_1 * x$  will fit well?
- Let's try one
- Looks bad, try another one.

## If $y$ is a discrete variable

51 | 53

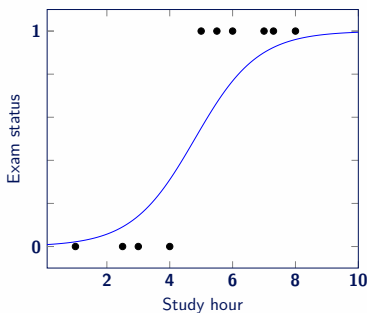


Exam status: 0 (fail), 1 (pass)

- Do you think the linear model  $y = \theta_0 + \theta_1 * x$  will fit well?
- Let's try one
- Looks bad, try another one.
- We see that such a linear model won't fit well, because of the nonlinear fashion. We need another model!

If  $y$  is a discrete variable

52 | 53



Exam status: 0 (fail), 1 (pass)

- We can use the logistic (sigmoid) function that reads

$$h_{\theta}(x) = g(\theta^T x) = \frac{1}{1 + e^{-\theta^T x}}$$

- We typically predict “1” if  $h_{\theta}(x) \geq 0.5$

1. Pick an initial line/model  $\mathbf{h}(\boldsymbol{\theta})$  by randomly choosing parameter  $\boldsymbol{\theta}$
2. Compute the corresponding cost function, e.g.,

$$J(\boldsymbol{\theta}) = \sum_{i=1}^m -\mathbf{y} \log \mathbf{h}_{\boldsymbol{\theta}}(\mathbf{x}^{(i)}) - (\mathbf{1} - \mathbf{y}) \log(\mathbf{1} - \mathbf{h}_{\boldsymbol{\theta}}(\mathbf{x}^{(i)}))$$

3. Update the line/model  $\mathbf{h}(\boldsymbol{\theta})$  by updating  $\boldsymbol{\theta}$  that makes  $J(\boldsymbol{\theta})$  smaller, using, e.g., the gradient descent that reads

$$\theta_j := \theta_j - \alpha \frac{\partial}{\partial \theta_j} J(\boldsymbol{\theta}),$$

where  $\alpha$  is the learning rate.

4. Repeat steps 2 and 3 until converges