## Learning Simple Functions II

Fundamentals of Artificial Intelligence

Fabien Cromieres

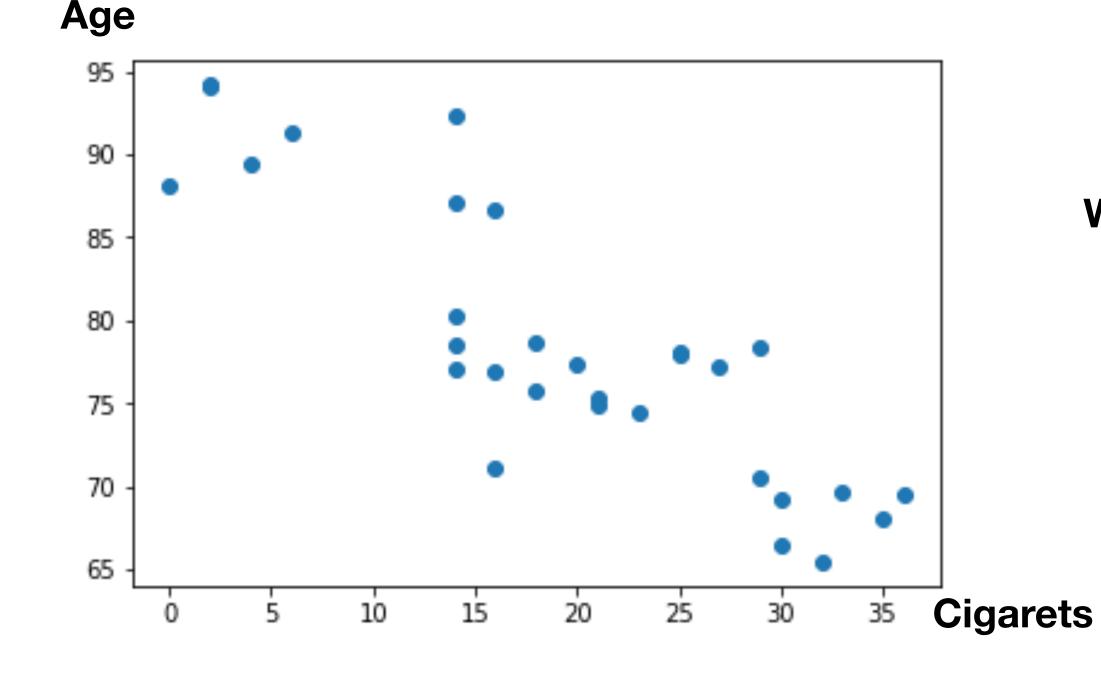
**Kyoto University** 

http://lotus.kuee.kyoto-u.ac.jp/~fabien/lectures/IA/

## Summary of Last Session 1/3

- Last session, we considered a simple case of Linear Regression:
- For 30 people, we know how old they died, and how many cigarets they were smoking per day
- Using this data, we wanted to predict the life expectancy of somebody given how many cigarets they were smoking

daily cigarets	age of death
32.0	73.471399
7.0	88.237207
30.0	82.077261
17.0	85.576741
27.0	76.190373
15.0	84.899030
20.0	72.598501
28.0	77.018773



What age am I most likely to die if I smoke 10 cigarets per days?

## Summary of Last Session 2/3

- We suppose there was some proportionality between the number of cigarets smoked and the reduction in life expectancy
- We suppose a linear relation between age of death and number of cigarets smoked

**Cigarets** 

$$age = \theta_0 + \theta_1 \times cig$$
Age
$$\frac{95}{90}$$
85
$$\frac{80}{75}$$
70
$$\frac{65}{90}$$
91
$$\frac{15}{20}$$
25
30
35
40

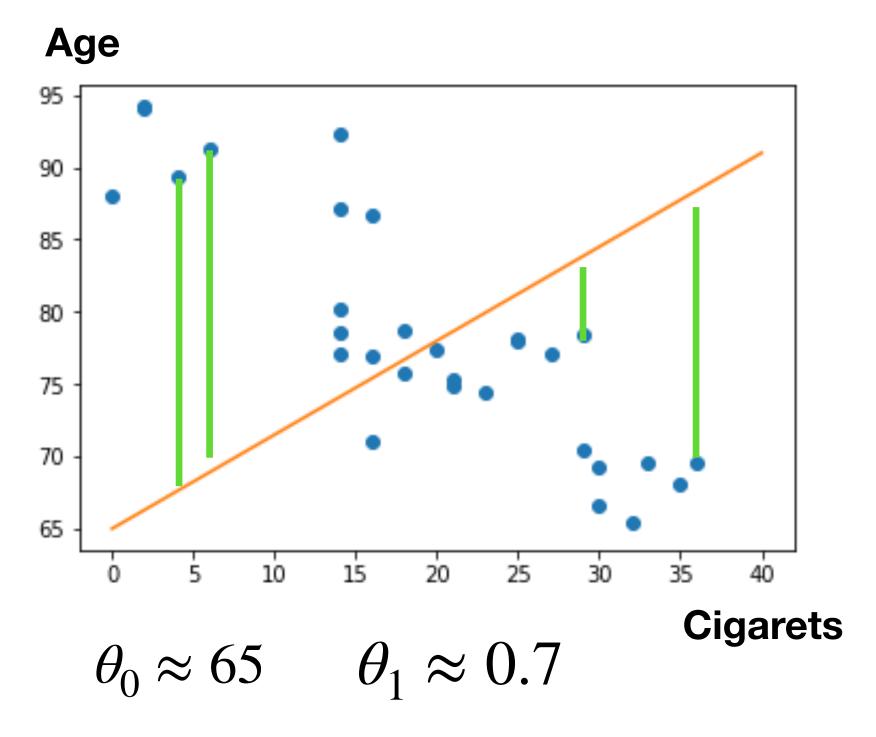
$$\theta_0 \approx 90$$

$$\theta_1 \approx -0.7$$

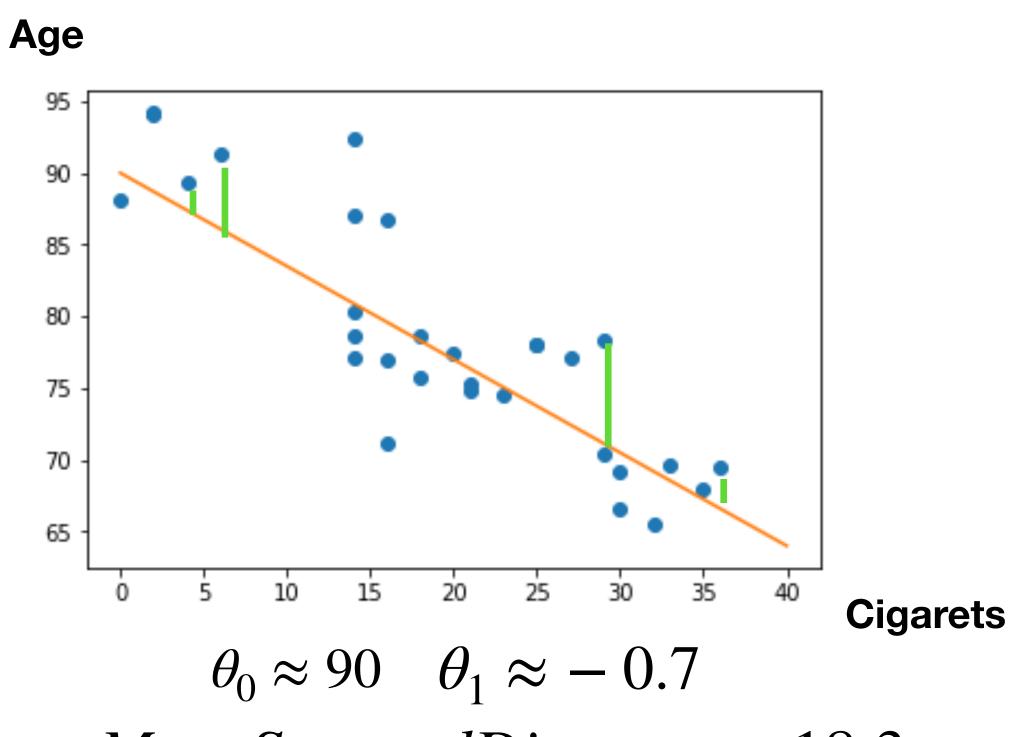
- Now, we can make our prediction:
- What age am I most likely to die if I smoke 10 cigarets per day?
  - 90-0.7x10 = 83 year old

## Summary of Last Session 3/3

- We saw that we could find the parameters of this linear relation by minimizing the Mean Squared Distance between the prediction of the model and the actual values
- We saw we could use Gradient Descent to do this minimization



Mean Squared Distance = 294.7



Mean Squared Distance = 18.3

## Today

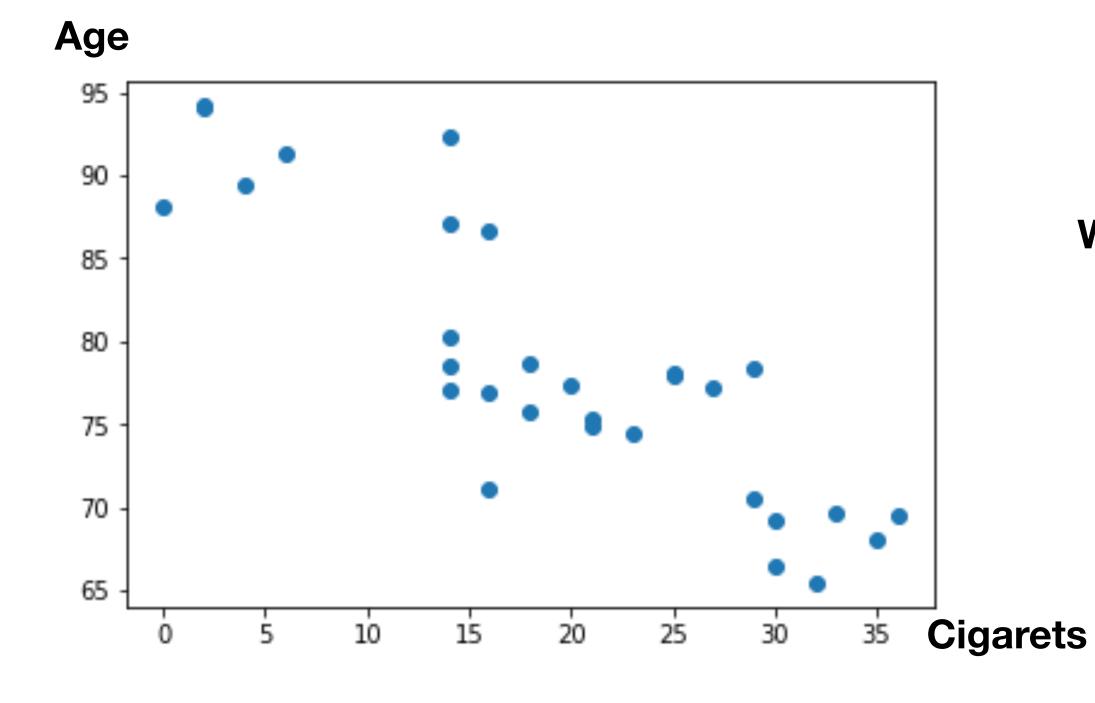
- We expand on this by considering:
  - More than one input feature
  - More complex functions
- We will have a look at the important concept of overfitting

## Adding more information

- The number of cigarets smoked is not the only important factor for predicting the age of death
  - Physical shape
  - Biological Sex
  - Wealth

•

daily cigarets	age of death
32.0	73.471399
7.0	88.237207
30.0	82.077261
17.0	85.576741
27.0	76.190373
15.0	84.899030
20.0	72.598501
28.0	77.018773



What age am I most likely to die if I smoke 10 cigarets per days?

## Adding more information

- Let us suppose now that, for the same 30 persons, we also have their BMI index and know their sex
- Vocabulary Note: In Machine Learning, we often cause a "<u>feature</u>" or "<u>feature function</u>" each of this piece of information about an example

	daily cigarets	bmi	is male	age of death
0	5.0	18.5	1.0	79.8
1	9.0	45.1	0.0	56.8
2	38.0	14.2	0.0	61.4
3	12.0	48.5	1.0	37.5
4	34.0	19.2	0.0	68.4
5	5.0	38.6	0.0	69.3
6	31.0	33.8	1.0	54.8
7	25.0	33.6	1.0	63.0
8	24.0	45.2	1.0	39.3
9	• • •	• • • •	• • • •	• • • •

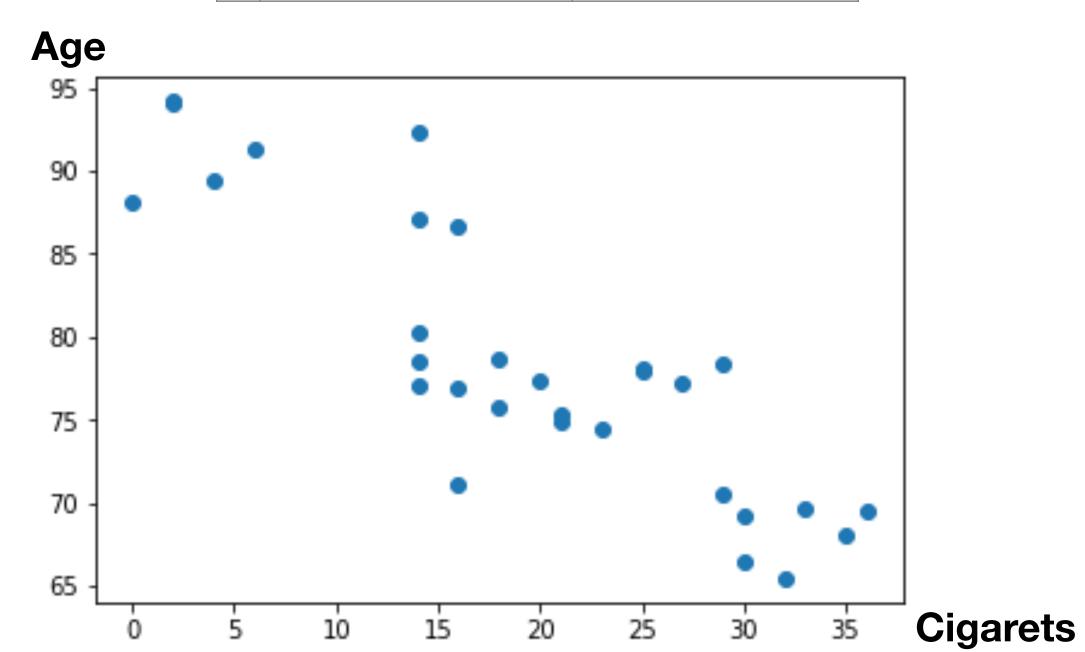
$$bmi = \frac{weight}{height^2}$$

Is male: 1 if the person is a male, 0 if female

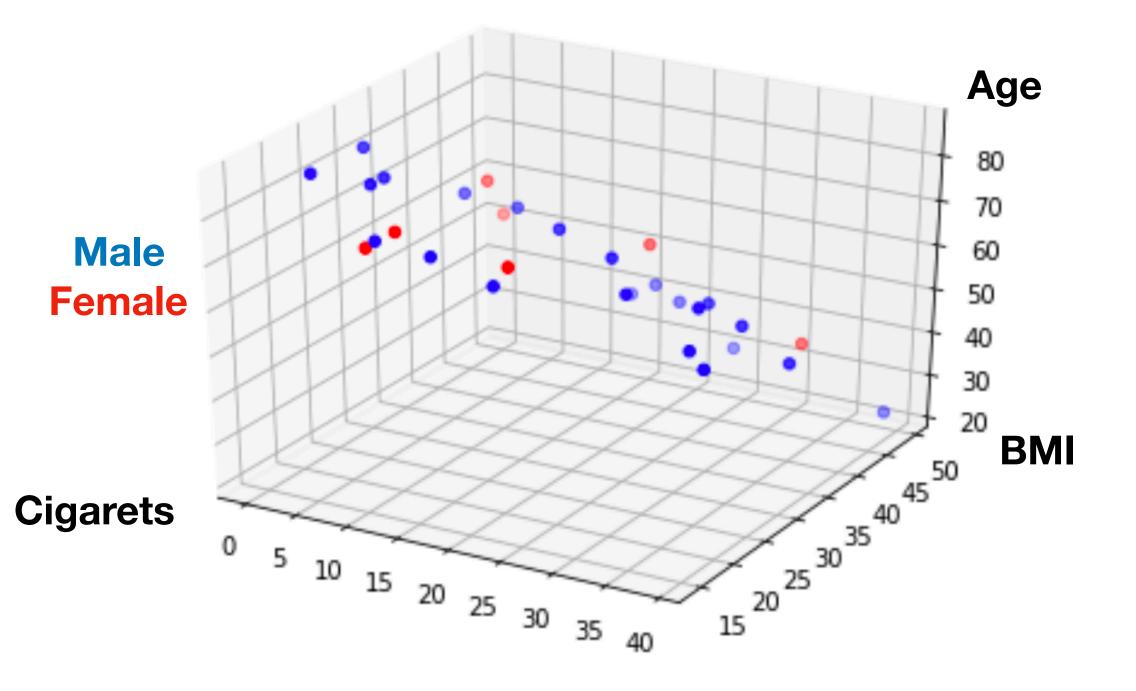
Note: most of the time, this is how we represent "categorical data" in Machine Learning: a feature equal to 1 if the example belongs to the category in question, and equal to zero otherwise

## Visual representation of the examples

	daily cigarets	age of death
0	5.0	79.8
1	9.0	56.8
2	38.0	61.4
3	12.0	37.5
7	25.0	63.0
8	24.0	39.3
9	• • •	• • • •



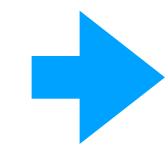
	daily cigarets	bmi	is male	age of death
0	5.0	18.5	1.0	79.8
1	9.0	45.1	0.0	56.8
2	38.0	14.2	0.0	61.4
3	12.0	48.5	1.0	37.5
7	25.0	33.6	1.0	63.0
8	24.0	45.2	1.0	39.3
9	• • •	• • • •	• • • •	• • • •



#### Linear Regression with more than one feature

 We can still suppose that the age of death is a linear function of the features

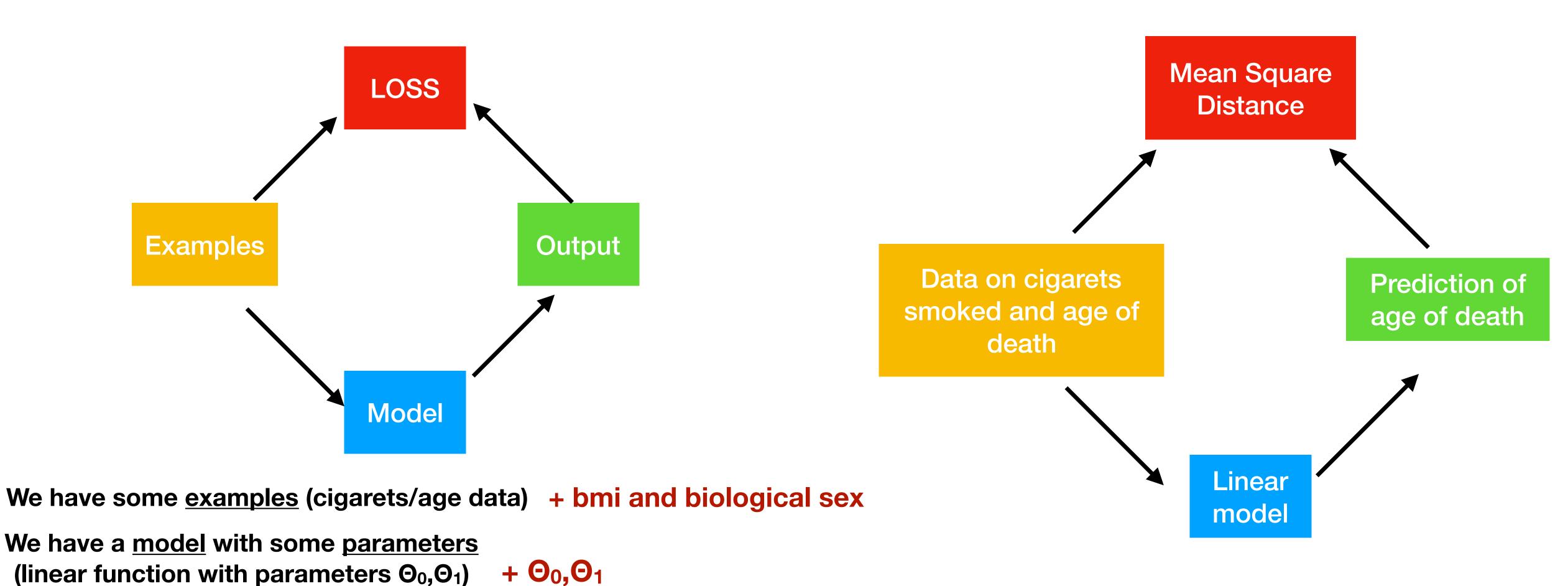
$$age = \theta_0 + \theta_1 \times cig$$



$$age = \theta_0 + \theta_1 \times cig + \theta_2 \times bmi + \theta_3 \times ismale$$

- We now have 2 more parameters (because we have 2 more features)
- For a total of <u>4 parameters</u>
- But finding the parameters will be done in exactly the same way as in the case with one feature

#### The General Supervised Machine Learning setting:

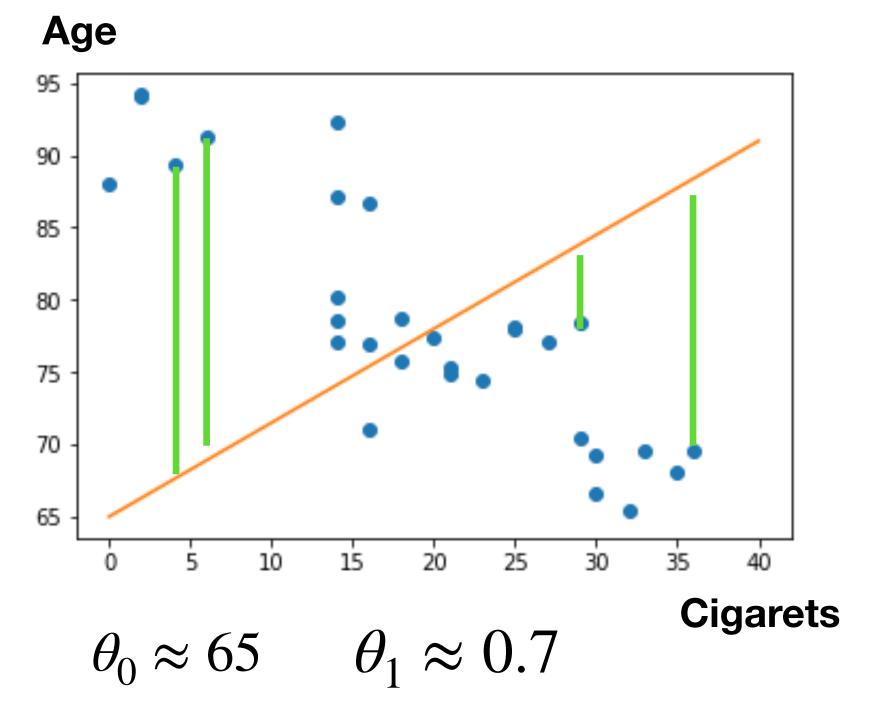


We have a <u>loss function</u> that compute the difference between the example and the prediction of our model (Mean Squared Distance)

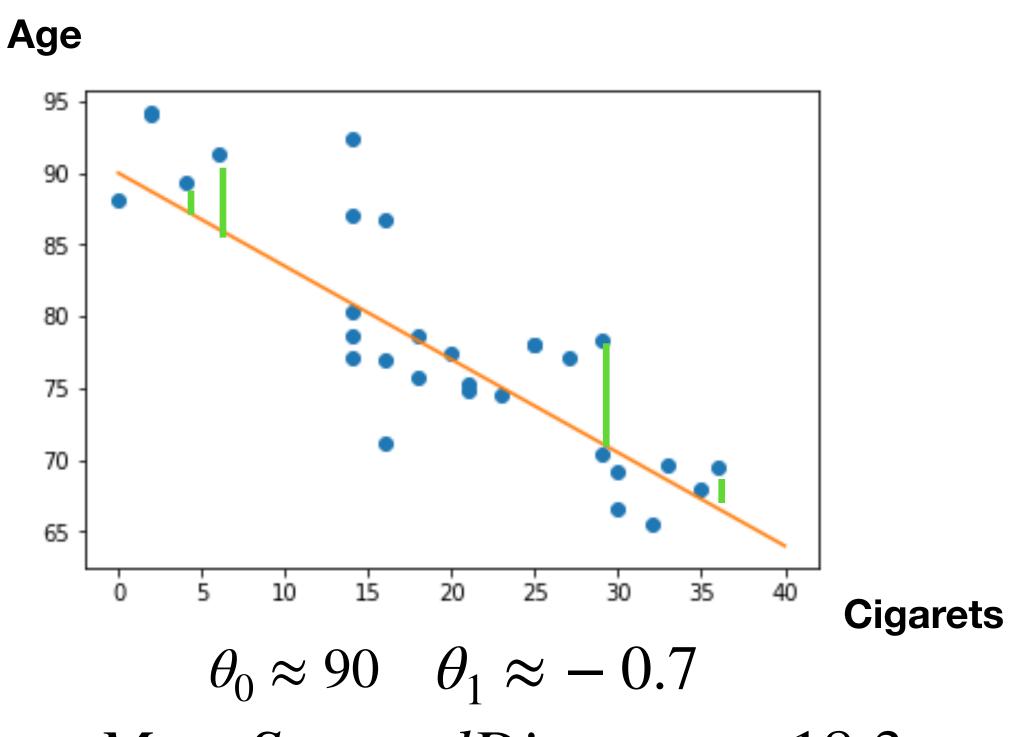
We minimize the loss to obtain the best parameters for our model

## Reminder: loss with one feature

- We saw that we could find the parameters of this linear relation by minimizing the Mean Squared Distance between the prediction of the model and the actual values
- We saw we could use Gradient Descent to do this minimization



Mean Squared Distance = 294.7



Mean Squared Distance = 18.3

#### Linear Regression with more than one feature

 The loss is defined as before: the average of the squared distance between model prediction and real example values

$$Mean Squared Distance = \frac{1}{N} \cdot \sum_{i} (model(cig_i, bmi_i, ismale_i) - age_i)^2$$
 Model prediction for example i. Real value for example i.

	daily cigarets	bmi	is male	age of death
0	5.0	18.5	1.0	79.8
1	9.0	45.1	0.0	56.8
2	38.0	14.2	0.0	61.4
3	12.0	48.5	1.0	37.5
7	25.0	33.6	1.0	63.0
8	24.0	45.2	1.0	39.3
9	• • •	• • •	• • • •	• • • •

$$age_{model} = \theta_0 + \theta_1 \times cig + \theta_2 \times bmi + \theta_3 \times ismale$$

## Small Exercise: compute the loss

 The loss is defined as before: the average of the squared distance between model prediction and real example values

$$Mean Squared Distance = \frac{1}{N} \cdot \sum_{i} \underbrace{(model(cig_i, bmi_i, ismale_i) - age_i)^2}_{\text{Model prediction for example i.}} - age_i)^2$$

	daily cigarets	bmi is	male	age of death	agraphian	$=\theta_0+\theta_1\times$	$cig + \theta$	$\sim \times hmi$	$+\theta_3 \times ismale$
0	5.0	18	0.0	80	"8° model	$-00 \cdot 01$	ceg i o	2 / 01100	1 03 / 15/11/01/0
1		45		57	$\alpha$	Λ 1		Λ 1	
2	38.0	16	0.0	61	$\theta_0 = 90$	$\theta_1 = -1$	$\theta_2 =$	<b>–</b> 0.1	$\theta_3 = -6$

Compute the model prediction for each of the examples.

Then, compute the loss.

$$Mean Squared Distance = \frac{1}{N} \cdot \sum_{i} \underbrace{(model(cig_i, bmi_i, ismale_i) - age_i)^2}_{\text{Model prediction for example i.}} \\ \text{Real value for example i.}$$

	daily cigarets	bmi	is male	age of death
0	5.0	18	0.0	80
1	9.0	45	1.0	57
2	38.0	16	0.0	61

$$age_{model} = \theta_0 + \theta_1 \times cig + \theta_2 \times bmi + \theta_3 \times ismale$$

$$\theta_0 = 90$$
  $\theta_1 = -1$   $\theta_2 = -0.1$   $\theta_3 = -6$ 

# What is the gradient?

 To do it a bit differently than last time: let us note the error on example i as:

$$error_i = model(cig_i, bmi_i, ismale_i) - age_i$$

• Then our loss is equal to: 
$$MeanSquaredDistance = \frac{1}{N} \cdot \sum_{i} (error_{i})^{2}$$

Then our loss is equal to:

$$\frac{\partial}{\partial \theta_k} Mean Squared Distance = \frac{1}{N} \cdot \sum_{i} 2 \times error_i \times \frac{\partial}{\partial \theta_k} error_i$$

(The mean squared distance is actually also often called the mean squared error)

Using linearity and the fact that

$$\frac{d}{dx}[f(x)]^2 = 2 \times f(x) \times \frac{d}{dx}f(x)$$

# What is the gradient?

$$\frac{\partial}{\partial \theta_k} Mean Squared Distance = \frac{1}{N} \cdot \sum_{i} 2 \times error_i \times \frac{\partial}{\partial \theta_k} error_i$$

 $error_i = model(cig_i, bmi_i, ismale_i) - age_i$ 

$$age_{model} = \theta_0 + \theta_1 \times cig + \theta_2 \times bmi + \theta_3 \times ismale$$

$$\frac{\partial}{\partial \theta_0} error_i = 1$$

$$\frac{\partial}{\partial \theta_0} error_i = cig_i$$

$$\frac{\partial}{\partial \theta_1} error_i = bmi_i$$

$$\frac{\partial}{\partial \theta_2} error_i = ismale_i$$

(we know have 4 parameters, so gradient is a 4-dimesnional vector)

$$gradient = \begin{bmatrix} \frac{1}{N} \cdot \sum_{i} 2 \times error_{i} \times 1 \\ \frac{1}{N} \cdot \sum_{i} 2 \times error_{i} \times cig_{i} \\ \frac{1}{N} \cdot \sum_{i} 2 \times error_{i} \times bmi_{i} \\ \frac{1}{N} \cdot \sum_{i} 2 \times error_{i} \times bmi_{i} \\ \frac{1}{N} \cdot \sum_{i} 2 \times error_{i} \times ismale_{i} \end{bmatrix}$$

## Small Exercise: compute the gradient

 $error_i = model(cig_i, bmi_i, ismale_i) - age_i$   $age_{model} = \theta_0 + \theta_1 \times cig + \theta_2 \times bmi + \theta_3 \times ismale$ 

• Let us compute the gradient for this examples and  $\Theta_k$ :

	daily cigarets	bmi	is male	age of death
0	5.0	18	0.0	80
1	9.0	45	1.0	57
2	38.0	16	0.0	61

Note you already computed the errors when you computed the loss in previous exercise

$$\theta_0 = 90$$
  $\theta_1 = -1$   $\theta_2 = -0.1$   $\theta_3 = -6$ 

$$gradient = \begin{bmatrix} \frac{1}{N} \cdot \sum_{i} 2 \times error_{i} \times 1 \\ \frac{1}{N} \cdot \sum_{i} 2 \times error_{i} \times cig_{i} \\ \frac{1}{N} \cdot \sum_{i} 2 \times error_{i} \times bmi_{i} \\ \frac{1}{N} \cdot \sum_{i} 2 \times error_{i} \times ismale_{i} \end{bmatrix}$$

Let us compute the gradient for this examples and Ok.

	daily cigarets	bmi	is male	age of death
0	5.0	18	0.0	80
1	9.0	45	1.0	57
2	38.0	16	0.0	61

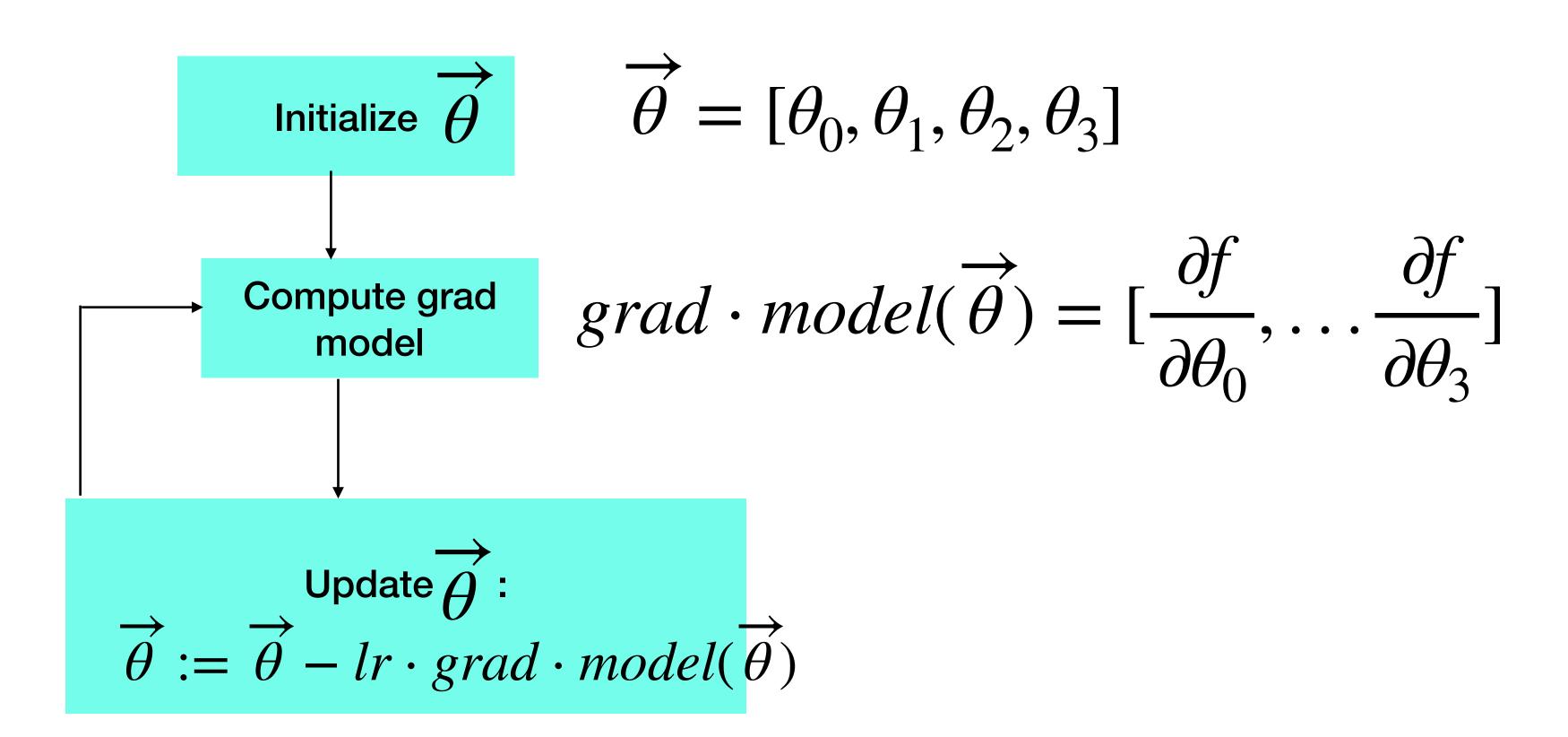
Note you already computed the errors when you computed the loss in previous exercise

$$\theta_0 = 90$$
  $\theta_1 = -1$   $\theta_2 = -0.1$   $\theta_3 = -6$ 

$$gradient = \begin{bmatrix} \frac{1}{N} \cdot \sum_{i} 2 \times error_{i} \times 1 \\ \frac{1}{N} \cdot \sum_{i} 2 \times error_{i} \times cig_{i} \\ \frac{1}{N} \cdot \sum_{i} 2 \times error_{i} \times bmi_{i} \\ \frac{1}{N} \cdot \sum_{i} 2 \times error_{i} \times ismale_{i} \end{bmatrix}$$

#### Gradient Descent

 Now that we know how to compute the gradient, we can find the optimal parameters with gradient descent:



## Feature Scaling

- We can see that our features have very different ranges and means:
  - "daily cigarets" is from 0 to 50
  - "bmi" is from 15 to 50
  - "is\_male" is from 0 to 1
- Usually, features having different ranges makes learning/gradient descent more difficult

	daily cigarets	bmi	is male	age of death
0	5.0	18.5	1.0	79.8
1	9.0	45.1	0.0	56.8
2	38.0	14.2	0.0	61.4
3	12.0	48.5	1.0	37.5
4	34.0	19.2	0.0	68.4
5	5.0	38.6	0.0	69.3
6	31.0	33.8	1.0	54.8
7	25.0	33.6	1.0	63.0
8	24.0	45.2	1.0	39.3
9	• • •	• • • •	• • • •	• • • •

The solution is to scale all the features to the same range

$$new = \frac{old - mean(old)}{max(old) - min(old)}$$

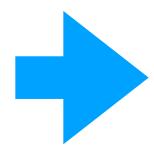
This way, all features have ranges between -1 and 1, and mean equal to zero

## Feature Scaling

After the features have been scaled, we apply Gradient Descent as usual:

 $new = \frac{old - mean(old)}{max(old) - min(old)}$ 

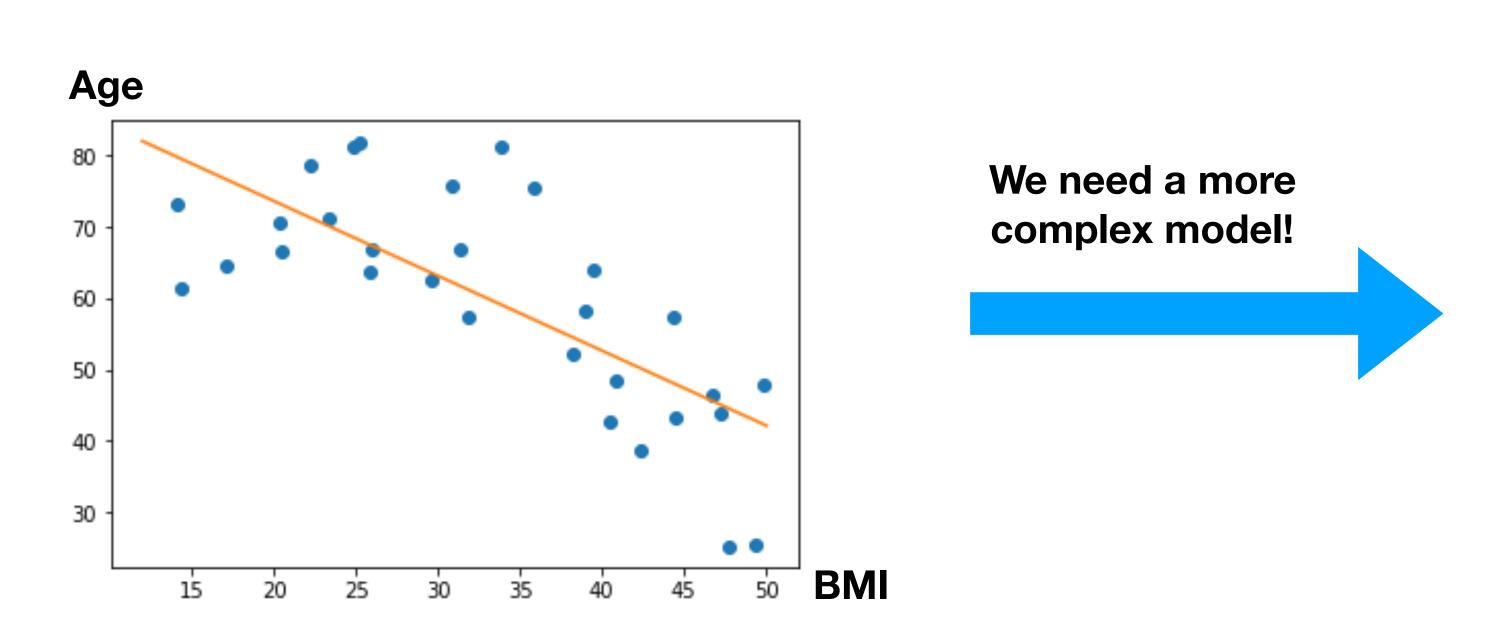
	daily cigarets	bmi	is male	age of death
0	24.0	44.4	1.0	43.4
1	37.0	47.7	1.0	25.2
2	11.0	14.1	0.0	73.1
3	33.0	49.3	1.0	25.5
4	27.0	20.4	0.0	70.5
5	27.0	38.2	1.0	52.1
6	6.0	22.2	1.0	78.6
7	31.0	17.1	0.0	64.6
8	28.0	23.3	0.0	71.3
9	15.0	31.4	1.0	66.8

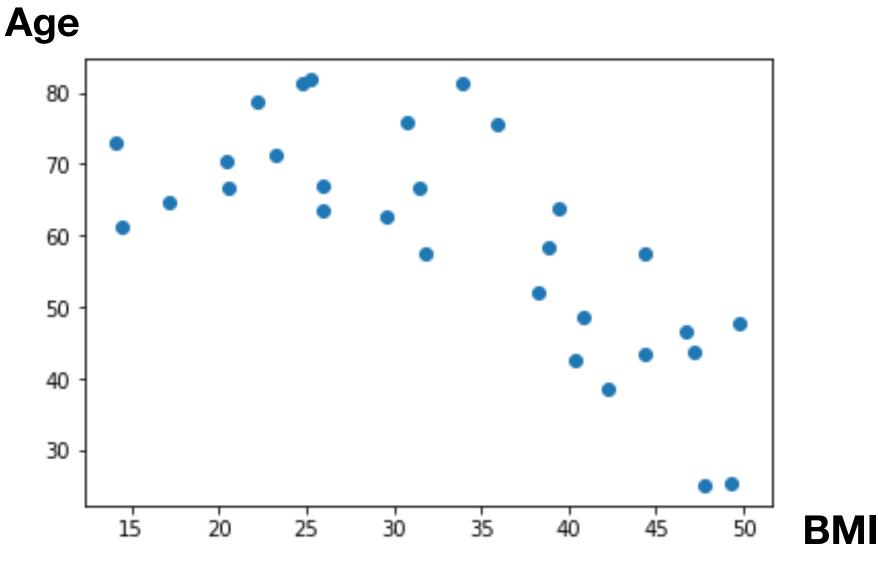


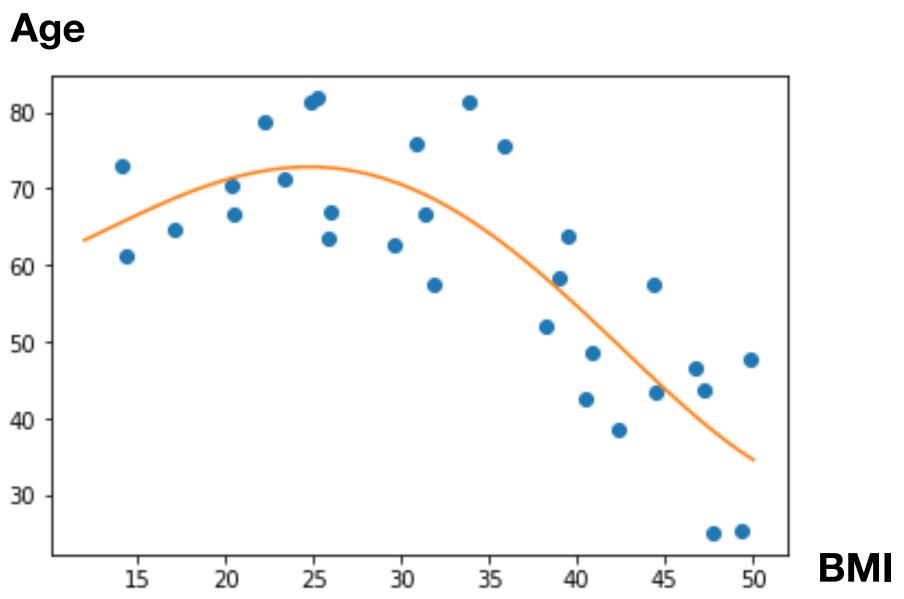
	daily cigarets	bmi	is male	age of death
0	0.064957	0.313072	0.433333	43.4
1	0.398291	0.405509	0.433333	25.2
2	-0.268376	-0.535668	-0.566667	73.1
3	0.295726	0.450327	0.433333	25.5
4	0.141880	-0.359197	-0.566667	70.5
5	0.141880	0.139402	0.433333	52.1
6	-0.396581	-0.308777	0.433333	78.6
7	0.244444	-0.451634	-0.566667	64.6
8	0.167521	-0.277965	-0.566667	71.3
9	-0.165812	-0.051074	0.433333	66.8

## Learning more complicated functions

- So far, we have only tried to learn linear functions of the data
- We might want to learn more complicated functions
- For example, having a high BMI reduce life expectancy
- But having a very low BMI also reduces life expectancy



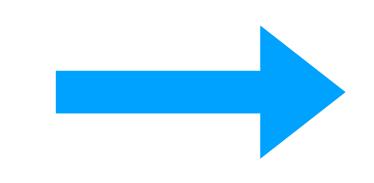




# Learning more complex functions: "Expanding the feature space"

- One neat trick to learn more complex functions: Create additional features from existing features
- Then apply linear regression
- Example: From the feature bmi, we add bmi<sup>2</sup>, bmi<sup>3</sup> and bmi<sup>4</sup>:

	bmi	age of death
0	44.4	43.4
1	47.7	25.2
2	14.1	73.1
3	49.3	25.5
4	20.4	70.5
5	38.2	52.1
6	•••	• • •



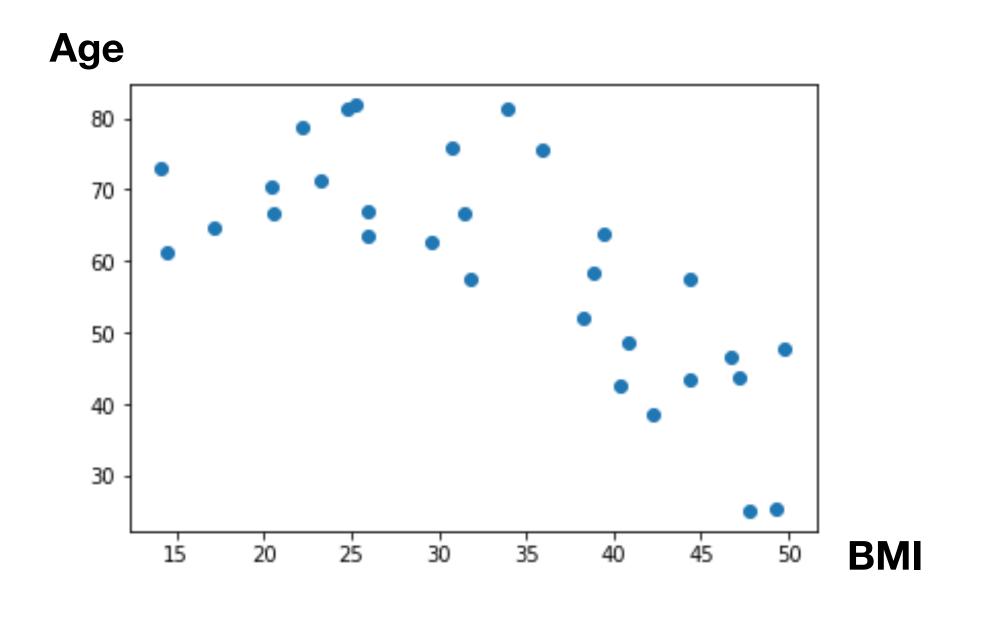
	bmi	bmi^2	bmi^3	bmi^4	age of death
0	44.4	1971.36	87528.384	3.88E+06	43.4
1	47.7	2275.29	108531.333	5.17E+06	25.2
2	14.1	198.81	2803.221	3.95E+04	73.1
3	49.3	2430.49	119823.157	5.90E+06	25.5
4	20.4	416.16	8489.664	1.73E+05	70.5
5	38.2	1459.24	55742.968	2.12E+06	52.1
6	• • •	• • •	• • •	• • •	• • •

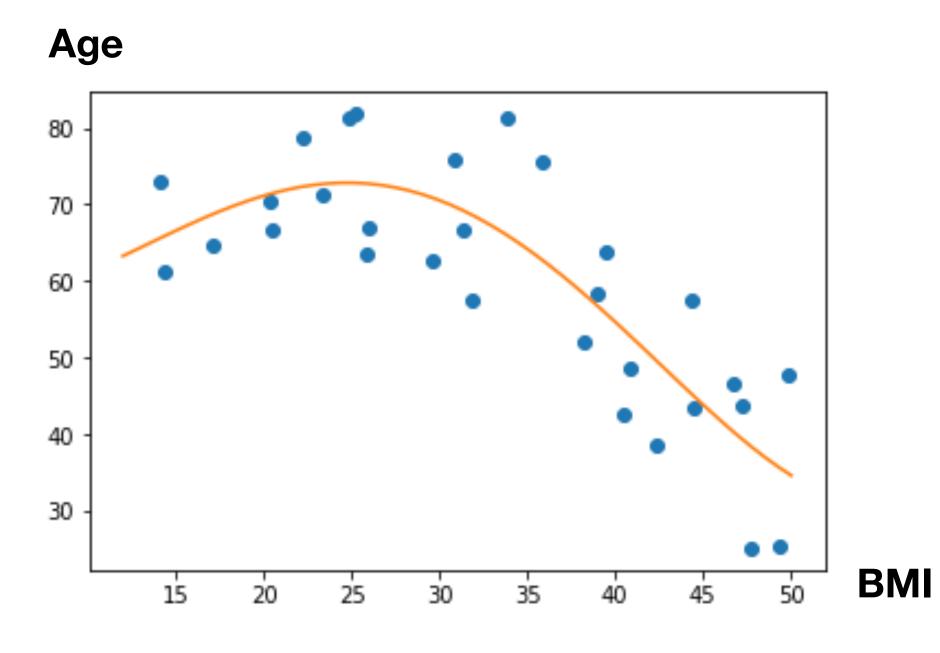
$$age_{model} = \theta_0 + \theta_1 \times bmi + \theta_2 \times bmi^2 + \theta_3 \times bmi^3 + \theta_4 \times bmi^4$$

## Learning more complicated functions

- This way, we can learn much more complex functions
- After finding the optimal parameters by gradient descent on the Mean Squared Distance:

$$age_{model} = \theta_0 + \theta_1 \times bmi + \theta_2 \times bmi^2 + \theta_3 \times bmi^3 + \theta_4 \times bmi^4$$





#### Learning Functions that are too complicated

Now, we might think:

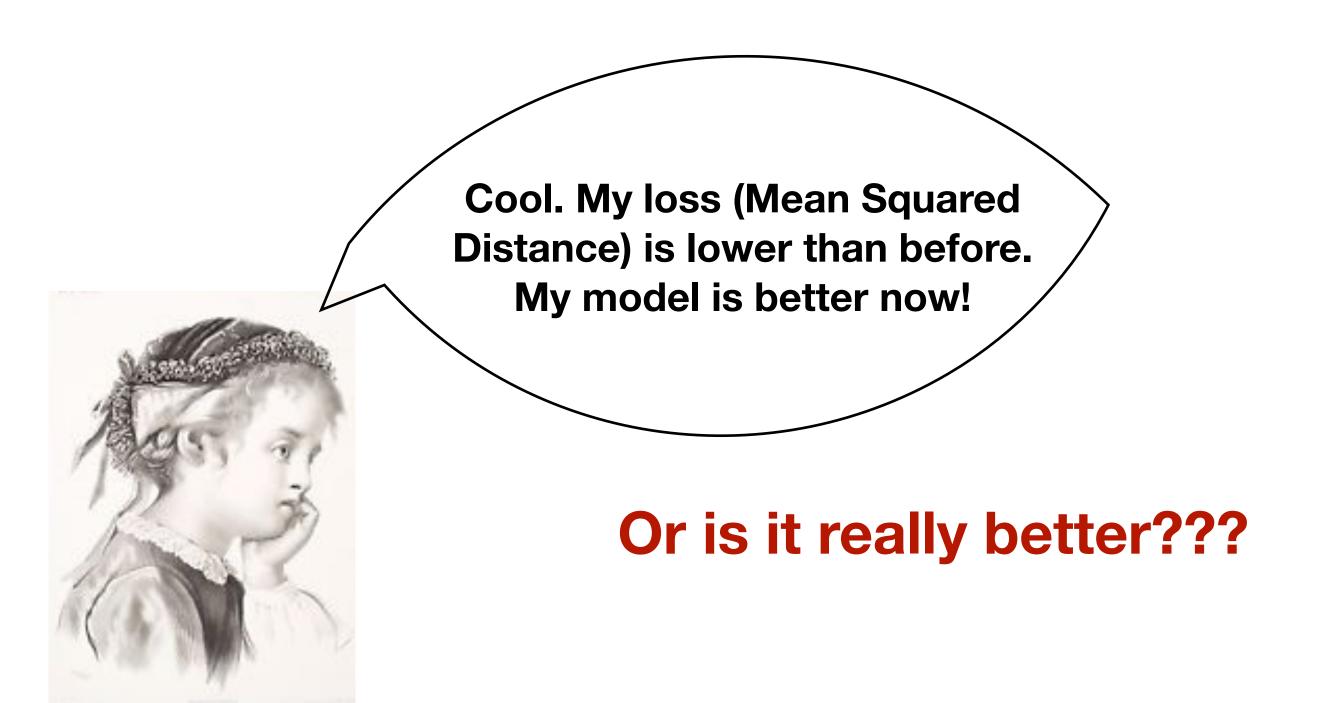
Cool. I am going to add as many variations of the features as I can. Then my loss (Mean Squared Distance) will drop to zero!

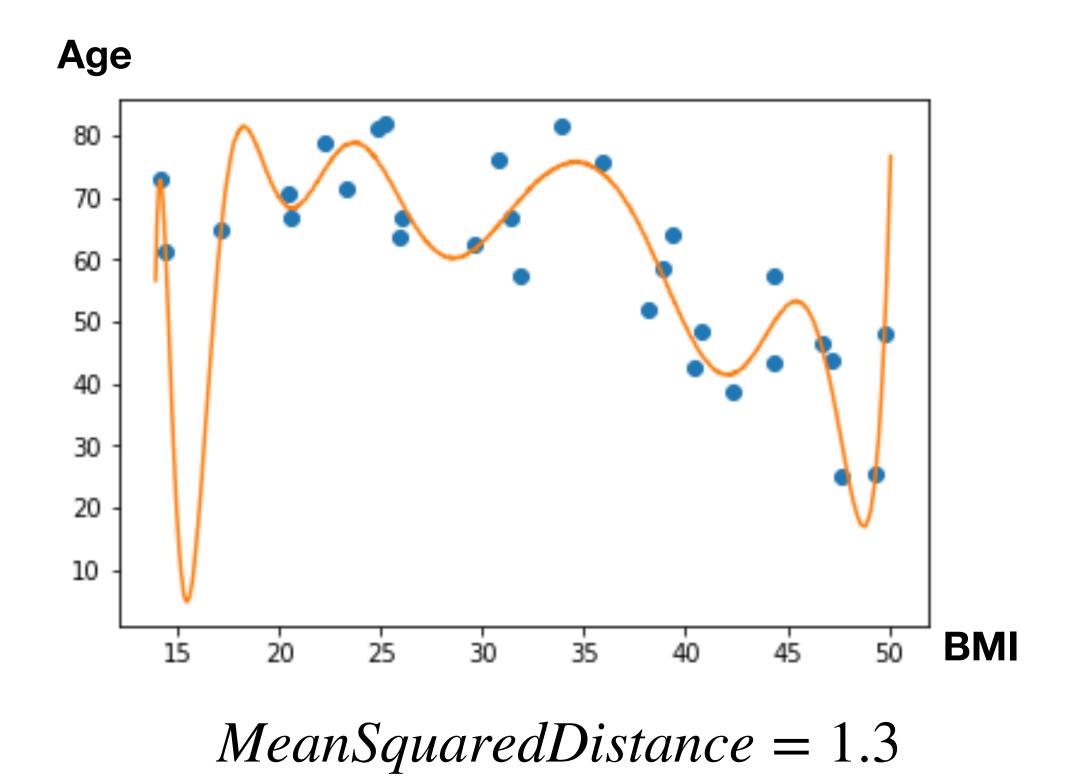
	bmi	bmi^2	bmi^3	bmi^4	bmi^5	sin(bmi)	log(bmi)	• • •	age of death
<b>0</b>	44.4	1971.36	87528.384	3.886260E+06	1.725500E+08	0.405662	3.793239	•••	43.4
1	47.7	2275.29	108531.333	5.176945E+06	2.469403E+08	-0.544766	3.864931	•••	25.2
2	14.1	198.81	2803.221	3.952542E+04	5.573084E+05	0.999309	2.646175	•••	73.1
3	49.3	2430.49	119823.157	5.907282E+06	2.912290E+08	-0.822324	3.897924	•••	25.5
4	20.4	416.16	8489.664	1.731891E+05	3.533059E+06	0.999793	3.015535	•••	70.5
5	38.2	1459.24	55742.968	2.129381E+06	8.134237E+07	0.480205	3.642836	•••	52.1
6	22.2	492.84	10941.048	2.428913E+05	5.392186E+06	-0.207336	3.100092	•••	78.6
7	17.1	292.41	5000.211	8.550361E+04	1.462112E+06	-0.984065	2.839078	•••	64.6
8	•••	•••	• • •	•••	• • •	• • •	•••	• • •	

 $age_{model} = \theta_0 + \theta_1 \times bmi + \theta_2 \times bmi^2 + \theta_3 \times bmi^3 + \theta_4 \times bmi^4 + \theta_5 \times bmi^5 + \theta_6 \times sin(bmi) + \theta_7 \times log(bmi) + \dots$ 

#### Learning Functions that are too complicated

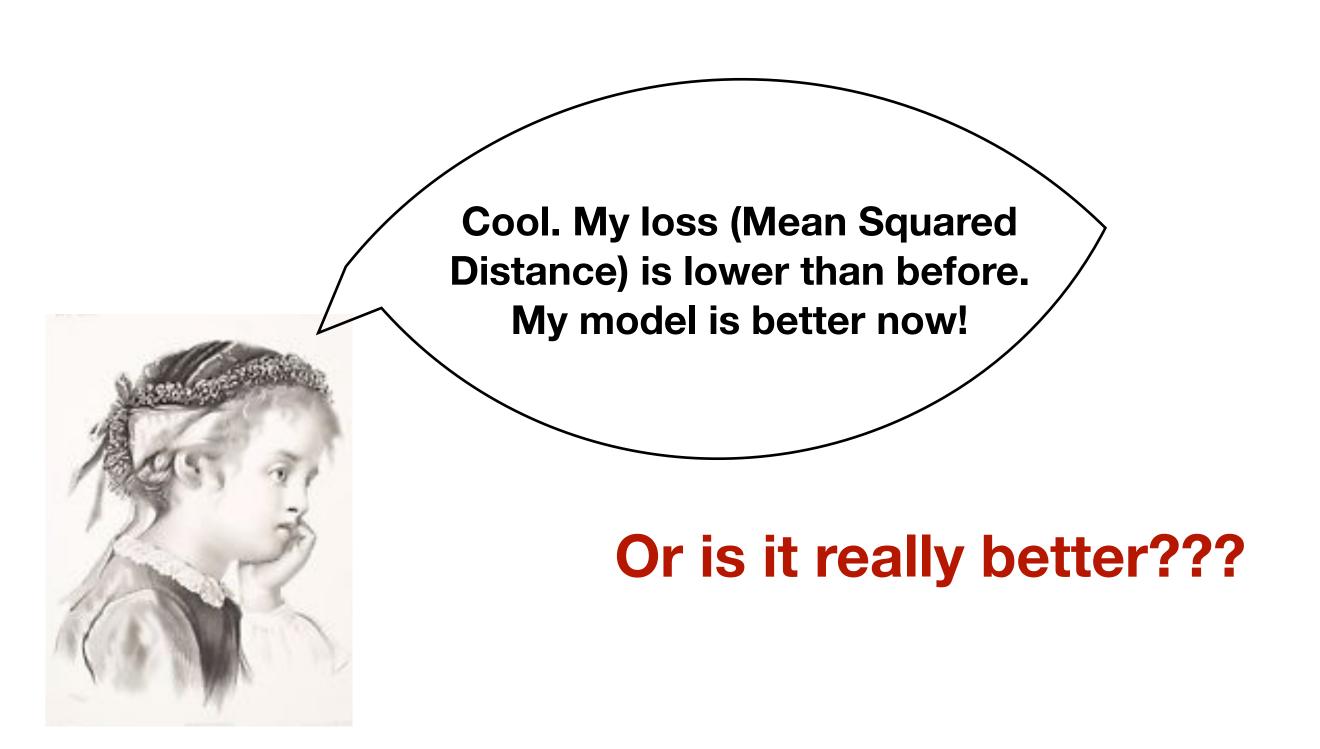
• The problem: this is what you are going to get:

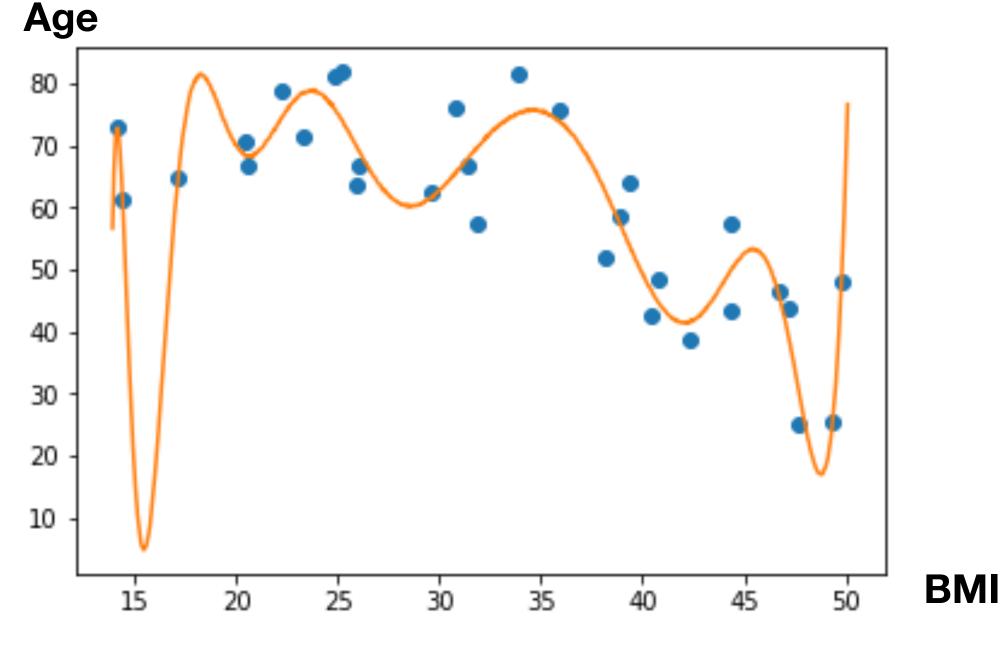




#### Learning Functions that are too complicated

• The problem: this is what you are going to get:



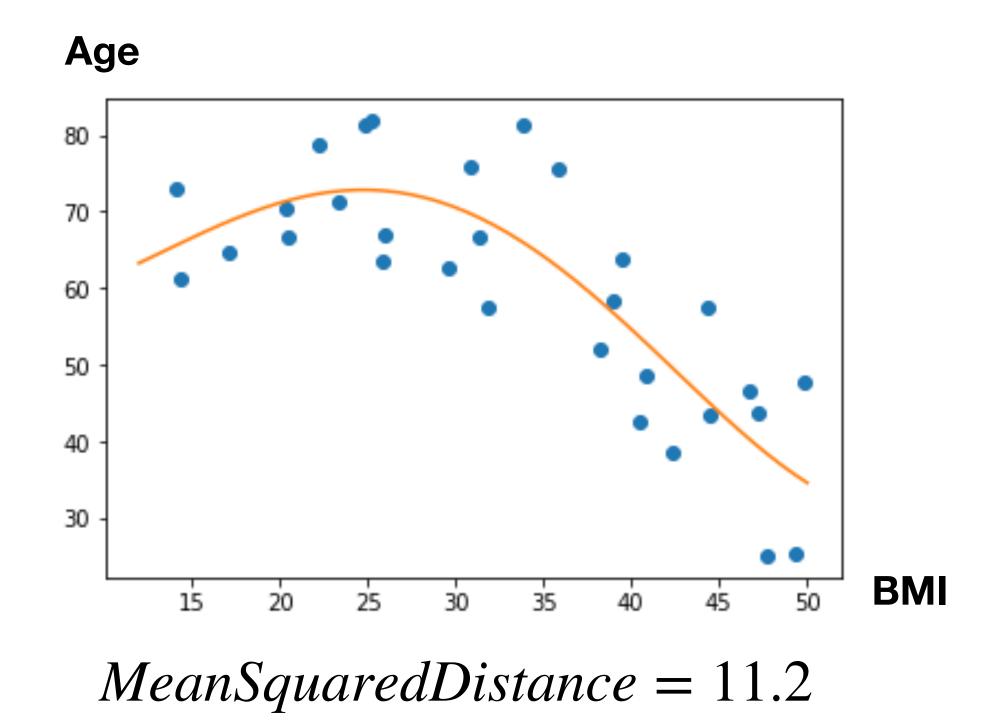


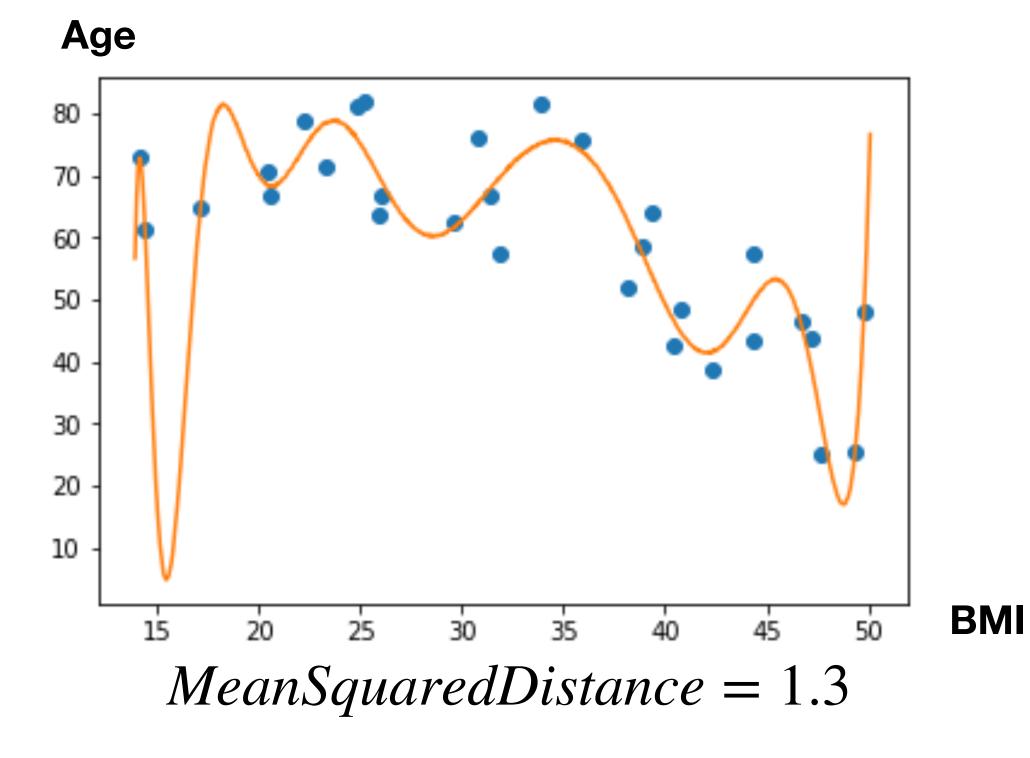
Mean Squared Distance = 1.3

The model predict that somebody with BMI 48 will die at 18 But somebody with BMI 51 is predicted to die at 80

## Overfitting

- It turns out that minimizing the loss does not always give the best model....
- This phenomenon is called overfitting





## Overfitting

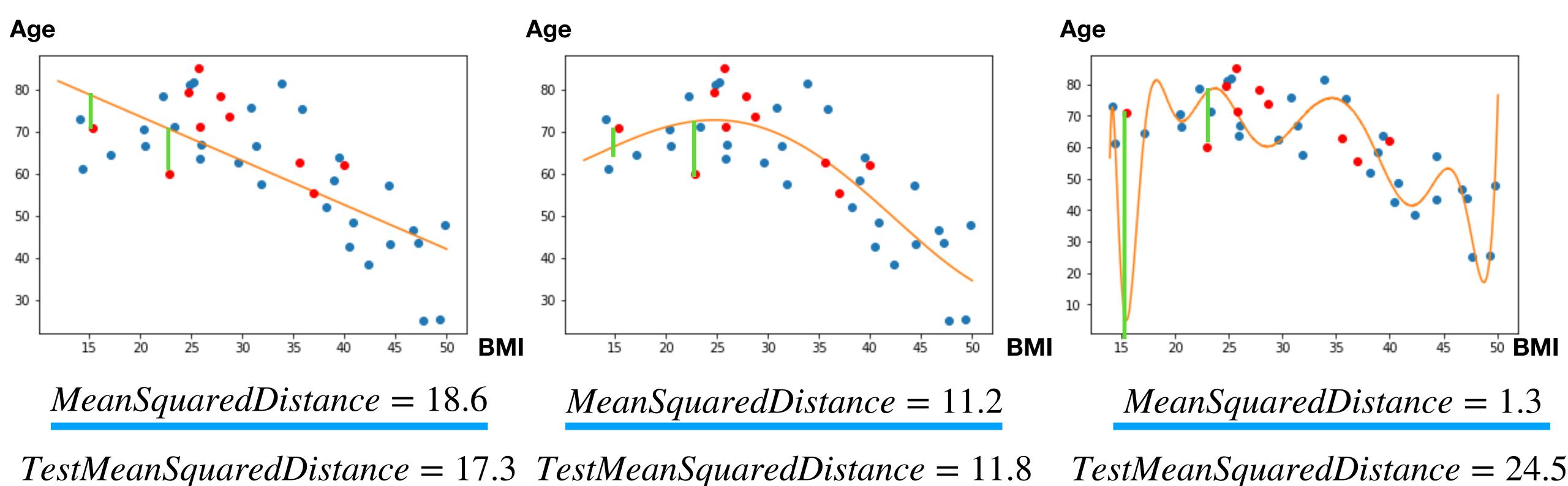
- Overfitting is to Machine Learning what Rote Learning is to Human Learning
- Instead of understanding the data, the model just memorized all of the examples
- If we ask it to predict the age of death of an example it has seen, it will give very good prediction
- But it will give very bad prediction as soon as we ask him to make a prediction for someone that was not in the example data.
- Very similar to a student that memorize without understanding the answer to a set of exercises in a class. He will do very well in the exam if the exam contains the exercises he studied, but very bad if the exercises are a bit different.

## How to detect overfitting?

- Pretty much like we do with humans: We evaluate them with different exercises than the ones they were trained for.
- In practice, it means we separate our example data in 2:
  - Training Data
  - Test Data
- We use the training data to do the Learning (we train the model)
- We use the test data to evaluate the quality of the learning (we test the model)

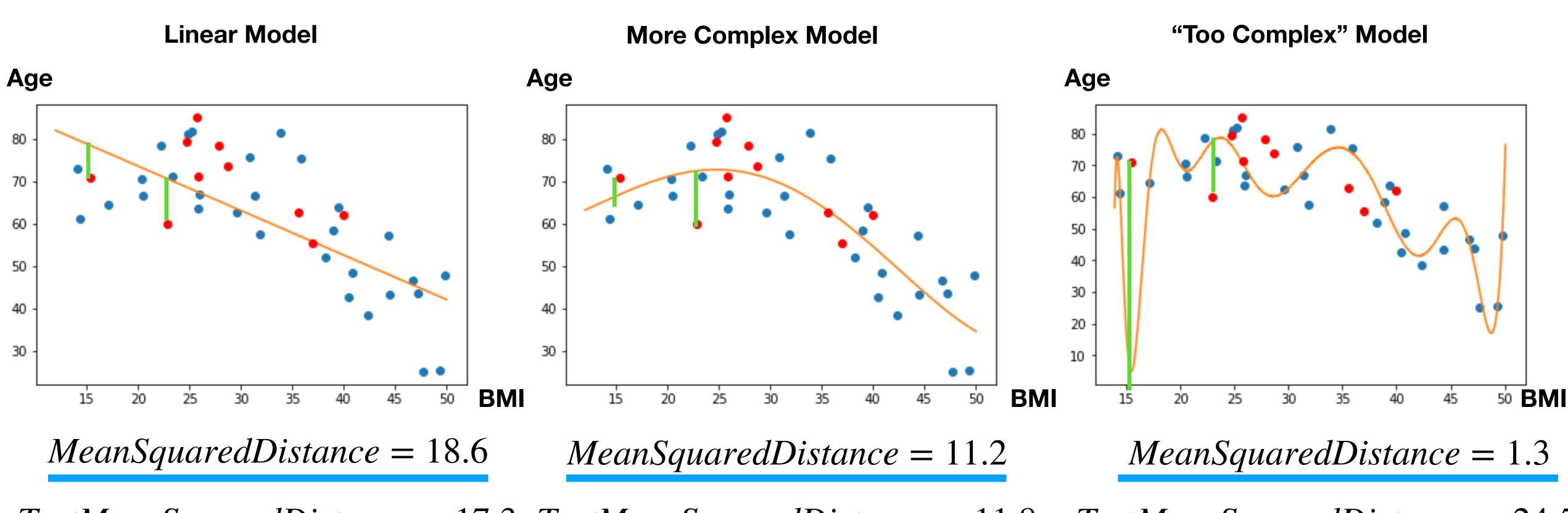
## How to detect overfitting?

- In blue: training examples
- In red: test examples



## How to detect overfitting?

 If we see a model with very low training loss and high test loss: it is overfitting!



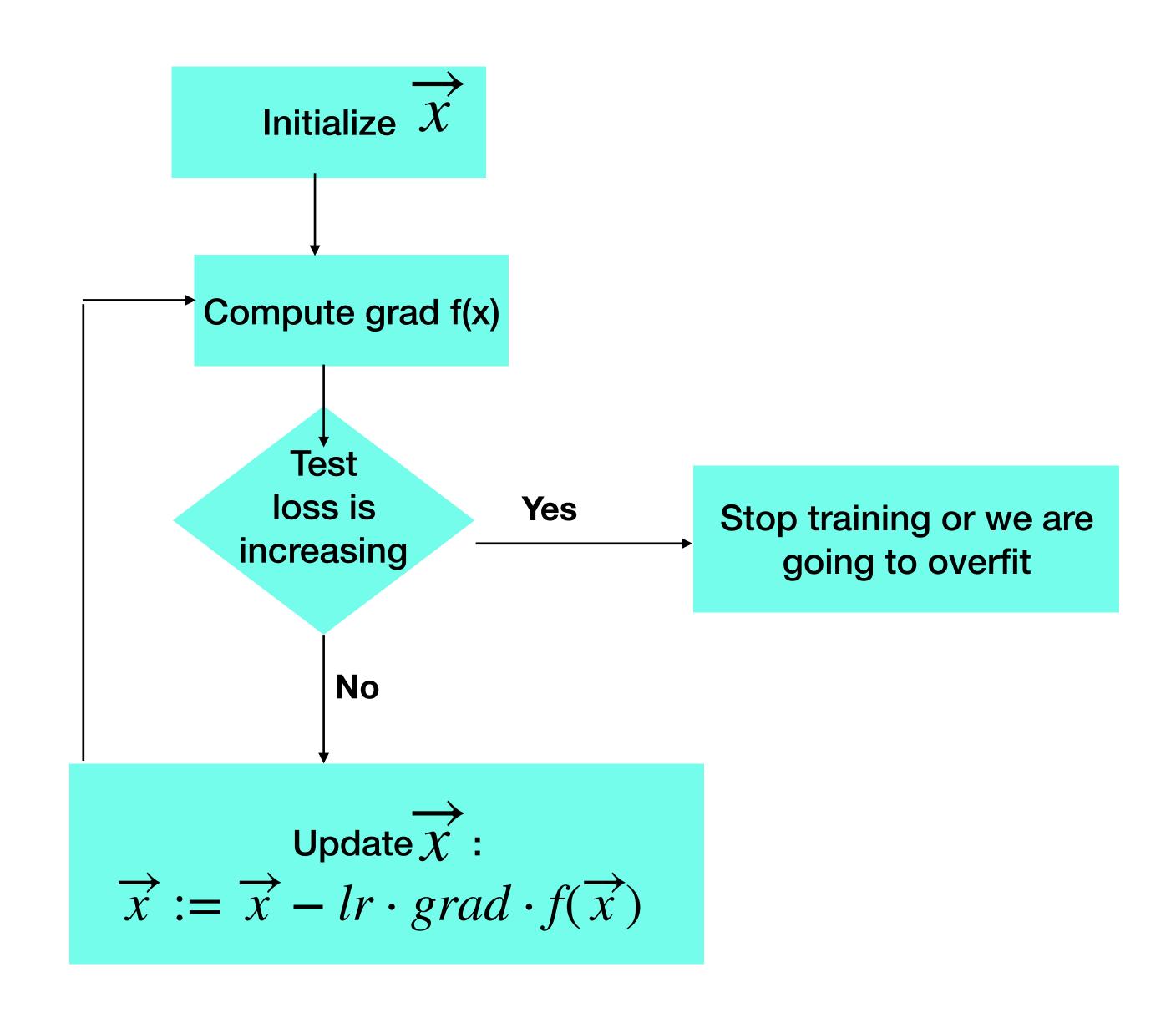
TestMeanSquaredDistance = 17.3 TestMeanSquaredDistance = 11.8

TestMeanSquaredDistance = 24.5

## How to prevent overfitting?

- A very simple and very efficient method for preventing overfitting is called "early stopping"
- During the gradient descent, <u>we check the test loss at each iteration</u>. When the <u>test loss starts to increase</u> (or stay unchanged) for a few iteration, we <u>stop</u> the training

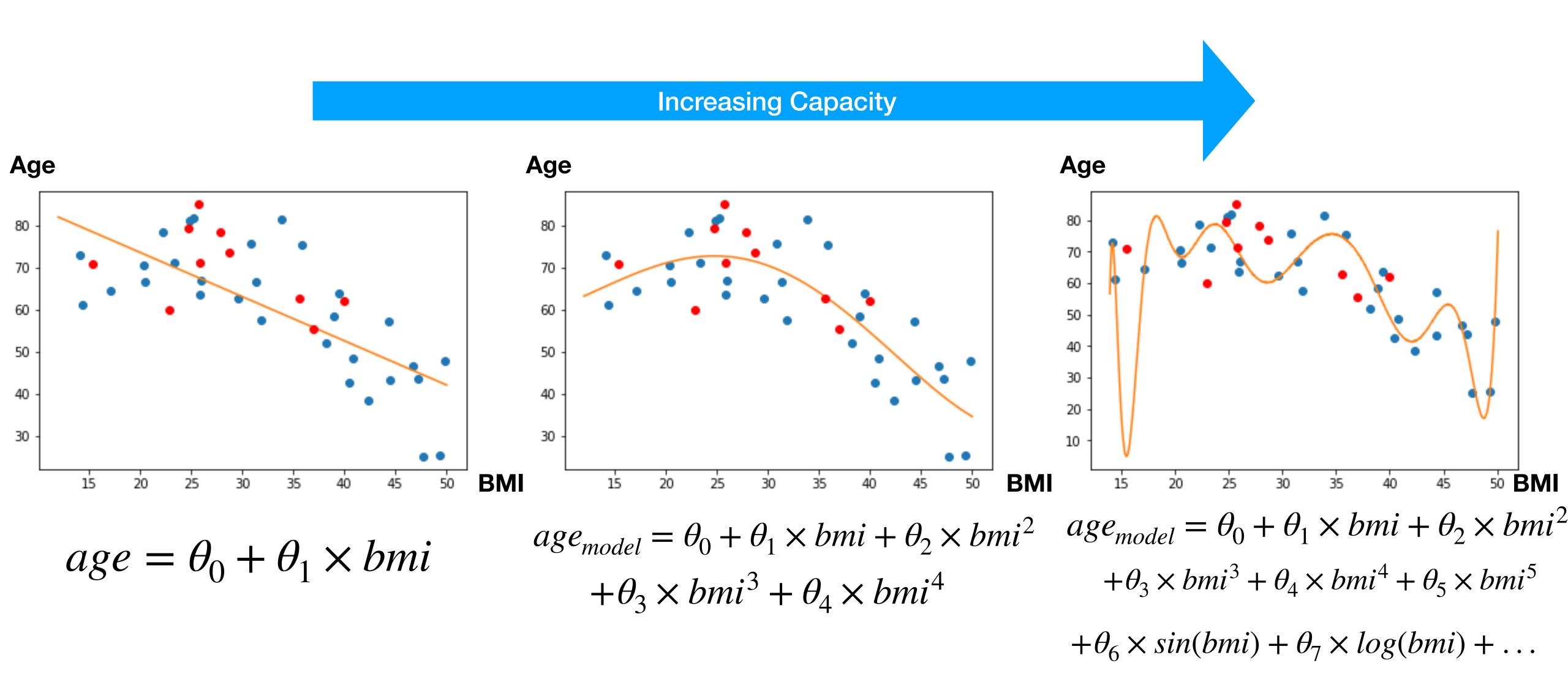
## Gradient Descent with Early Stopping



## How to prevent overfitting?

- Another method is to reduce the capacity of the model
- Roughly speaking, the capacity of a model is its <u>ability to adapt to a large</u> number of examples
- The capacity of a model will increase with the number of parameters

## How to prevent overfitting: Capacity



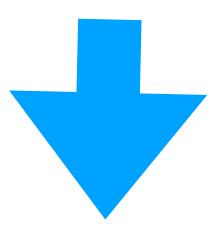
## How to prevent overfitting?

- Models with high capacity can learn more complicated relations
- But they overfit more easily
- A model with very high capacity has the ability to memorize all the training examples (the rote learning problem)
- If we reduce the capacity of a model, we can make it less prone to overfitting

## How to prevent overfitting: Capacity

 The simplest way to reduce capacity is to remove some parameters in the model:

$$age_{model} = \theta_0 + \theta_1 \times bmi + \theta_2 \times bmi^2 + \theta_3 \times bmi^3 + \theta_4 \times bmi^4 + \theta_5 \times bmi^5 + \theta_6 \times sin(bmi) + \theta_7 \times log(bmi) + \dots$$



$$age_{model} = \theta_0 + \theta_1 \times bmi + \theta_2 \times bmi^2 + \theta_3 \times bmi^4 + \theta_4 \times bmi^5$$

## How to prevent overfitting: Capacity

- Another way to reduce the capacity is "force" the model to use small values for the weights
  - By default, the parameters  $\Theta_k$  can take any value: -100 000, 0.1, 1 000 000
  - If we forbid the model to use very high values for  $\Theta_k$ , we reduce its capacity: it cannot adapt to data as well as before
- Most practical way to forbid high values: "L2 Regularization"

## L2 Regularization

- We note  $|\Theta|^2$  the sum of the square of all parameters  $\Theta_k$  (this is called the "L2 Norm")
- Then we add this quantity to the loss we want to minimize:

$$Loss = MeanSquaredDistance = \frac{1}{N} \cdot \sum_{i} (model(cig_{i}, bmi_{i}, ismale_{i}) - age_{i})^{2}$$
 
$$Loss = \frac{1}{N} \cdot \sum_{i} (model(cig_{i}, bmi_{i}, ismale_{i}) - age_{i})^{2} + \lambda |\overrightarrow{\theta}|^{2}$$

Then we apply Gradient Descent to this new loss

#### Next

- Next time, we will consider Classification Problems:
  - Predicting if some symptoms are the sign of a disease or not...
  - Predicting if an image represents a cat or a dog
  - Predicting if a text is in French, German or Japanese
- Also: a small assignment is coming this week (for a deadline in two weeks)
  - Will be simple. Just for a review of what we have seen so far.