

# Applied Machine Learning - Basic

Prof. Daniele Bonacorsi

## Lecture 5

Data Science and Computation PhD + Master in Bioinformatics  
**University of Bologna**

1110001110100100110011010010100111010100100100101011  
100101010010101010101101001010101011100011101001  
00110011010010100111010010010101110010100101001010  
11100110**Advices**00**for**01**applying**00**ML**10101010101110101  
01001010100110101001100100101011101010010010101010  
101101001010101011100011101001001100110100101001  
11010110110101110010101001010101010110001010101  
1110001110100100110011010010100111010100100101011  
10010101001001110101101101011100101001010101010  
11010001010111000111001001100100101001110101001  
0010101110010101001001110101101101011100101010010  
1010101101000101010111000111001001101001010010011  
10101001001011100101010010011101011011010111001  
0101001010101101000101011100011100101010101001000  
1101111000100110011010010100111010101011100010010  
11001010101010111010011011101010100111010111010111011

1110001110100100110011010010100111010100100100101011  
100101010010101010101101001010101011100011101001  
00110011010010100111010010010101110010100101001010  
111000 **Advices** 00 **for** 01 **applying** 00 **ML** 10101010101110101  
001010100110100110010010101110101010010101010  
110100101010101110001110100100110011001001010011  
010110110101110010100101010101011000101010111  
1000111010010011001101001010011101010010010101110  
0101010010011101011011011001010100101010101011  
01000101010 **Evaluating** 10 **a** 10 **learning** 01 **model** 00100100  
11010111001010100101010101101000101011100011  
10010100110 **Next** 10 **steps** 100110100011101001110101100  
1101011100101010010101010101101000101011100011  
100101010001101110001001100110100101001110101010  
11100010010110010101010101110011001101110101001  
11010111011010101010110010100011100001100011011

---

---

Knowing learning algos is not enough.

Application of ML on real-world datasets offer **no “gold recipes”**

- i.e. precise steps to follow that would solve your problem, whatever it is
- all efforts are (potentially) painful and time-consuming, and there is no obvious direction to follow..

Any ML practitioner (less or more expert) need to be able to realise quickly in which direction next efforts had better be invested.

What follows aims at giving a number of practical suggestions, advice, guidelines on how to improve in that.

Key question in this section is: **in developing a ML system and/or trying to improve the performances of a developed one, how do I decide what are the best next steps?**

---

---

Example: housing prices, again.

You have implemented and regularised linear regression, minimised the cost function  $J$ , tuned a bit your hyper-parameters.

But if you test your hypothesis on new set of houses, suppose you find that this is making huge errors in the prediction of the housing prices.

**What should you then try in order to improve the learning algo?**

# Improving performances of a learning algo

One thing you could try, is to **get more training examples**.

- really, stop trying to make the model work. Assess the amount of data you have, and decide. Think about how to get more.
  - ❖ “oh, if I had twice as many data...”

Sometimes getting more training data doesn't actually help. We will see.

Another things you might try is to **try a smaller set of features**.

- carefully spend some time to select a subset only, just to prevent overfitting with brute force.

Or maybe you need to **get additional (or different) features**

- maybe the current set of features aren't informative enough
- maybe you want to collect more data also in order to have a chance to extract more (or better) features

Sometimes a change in these directions requires a lot of work to get more data (e.g. surveys, research, etc..) so huge projects, and so once again it would be nice to know in advance if this is ultimately going to help before spending a lot of time pursuing this..

You can also try **adding polynomial features**

Try **decreasing/increasing lambda**

# So, what?!

---

Unfortunately, the most common method that people use to pick one of these is as simple as: **go by experience and gut feeling.**

Chances that all this turns out to be **very ineffective** are high..

# ML diagnostics

---

Be pragmatic.

*Isn't there a pretty simple technique that can let you very quickly select a fraction of the actions on the aforementioned list as being potentially promising actions to implement, and rule out all the rest?*

Fortunately yes!

We will talk about **machine learning diagnostics**, i.e. **how to evaluate learning algorithms**

- tests you can run, to get insight into what is or isn't working with an algorithm, which will often give you insight as of what are promising things to try to improve a learning algorithm's performance
- NOTE: diagnostics efforts can take time to implement and understand, but it is ultimately a good use of your time. Some ML, by the way, are more friendly than others to do so..

1110001110100100110011010010100111010100100100101011  
100101010010101010101101001010101011100011101001  
00110011010010100111010010010101110010101001010  
1110**Advices**00**for**01**applying**00**ML**1010101011101010100  
101010011010011001001010111010101001010101011  
010010101010111000111010010011001100100101001101  
0110110101110010101001010101011010001010101110  
0011101001001100110100101001110101001001010111001  
01010010011101011011010111001010100101010101101  
0111000010**Evaluating**10a10**learning**01**model**1001001001  
1010111001010100101010101011000101010111000111  
0010101110**Evaluating**10a00**hypothesis**10011010001110  
1001110101100110101110010101001010101010110100010  
10101111000111001010001101111000100110011010010  
1001110101011000100101100101010101010111010001  
110101010011101011101010101010110000111000010110

---

---

We have made a hypothesis.

This “has been learned” by my system (i.e. I have  $\theta$ s)

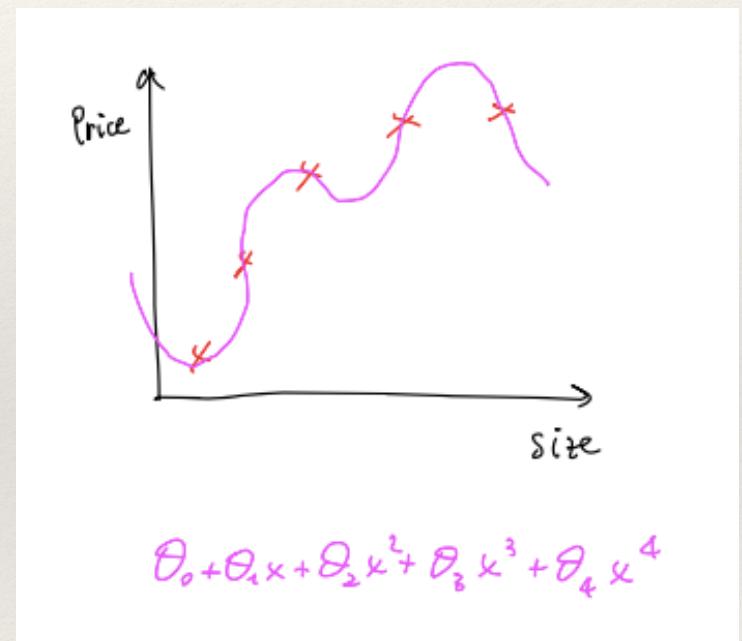
How to evaluate it?

# Check for overfitting

We fitted the parameters of our learning algorithm, and we extracted the parameters that minimise the training error.

Getting a really low value of training error for a hypothesis might be a good thing, but we have already seen that that doesn't mean it is a good hypothesis.

- E.g. overfitting -> it fails to generalise to new examples not in the training set



How do you tell if your hypothesis is suffering from overfitting?

- of course, plotting does not help in higher dimensions...

size (m <sup>2</sup> )	price (k€)
150	300
100	200
80	200
75	250
160	400
210	410
90	200
95	280
105	350
120	370

NOTE : you can have  
THOUSANDS ...

size (m <sup>2</sup> )	price (k€)		
150	300	70%	$(x^{(1)}, y^{(1)})$
100	200		:
80	200		:
75	250		
160	400	TRAINING SET	$(x^{(m)}, y^{(m)})$
210	410		$m = \# \text{ training examples}$
90	200		
95	280	TEST SET	$(x_{\text{test}}^{(1)}, y_{\text{test}}^{(1)})$
105	350		:
120	370	30%	$(x_{\text{test}}^{(m_{\text{test}})}, y_{\text{test}}^{(m_{\text{test}})})$

NOTE : if there is some ordering into the data,  
shuffle and after do the 70-30 split.

$m_{\text{test}} = \# \text{ test examples}$

# Training/testing procedure for linear regression

Procedure:

- learn parameters  $\theta$  from the training dataset (i.e. the 70% of the total)...
- ... and minimise the training error  $J(\theta)$  (or  $J_{\text{train}}(\theta)$  now, if you wish)
- then, compute the **test set error**

(linear regression (w/o regularization))

$$J_{\text{test}}(\theta) = \frac{1}{2m_{\text{test}}} \sum_{i=1}^{m_{\text{test}}} \left( h_{\theta}(x_{\text{test}}^{(i)}) - y_{\text{test}}^{(i)} \right)^2$$

# Training/testing procedure for logistic regression

Procedure is the same as before:

- learn parameters  $\theta$  from training dataset (i.e. the 70% of the total)...
- ... and minimise the training error  $J(\theta)$  (or  $J_{\text{train}}(\theta)$  now, if you wish)
- then, compute the **test set error**

logistic regression (w/o regularization)

$$J_{\text{test}}(\theta) = - \frac{1}{m_{\text{test}}} \sum_{i=1}^{m_{\text{test}}} y_{\text{test}}^{(i)} \log h_{\theta}(x_{\text{test}}^{(i)}) + (1 - y_{\text{test}}^{(i)}) \log h_{\theta}(x_{\text{test}}^{(i)})$$

.. OK, but not the best one ..



# 0/1 misclassification error

---

While this definition of the test set error  $J_{\text{test}}$  is perfectly reasonable, sometimes there is an alternative test set metric that might be easier to interpret:

the **0/1 misclassification error**, with the 0 and 1 labels denoting that you either get an example right or you get an example wrong.

0/1 misclassification error

a classification  
metrics

$$\underbrace{\text{err} (h_\theta(x), y)}_{\text{error of a prediction}} = \begin{cases} 1 & \begin{array}{l} \text{if } h_\theta(x) \geq 0.5, y=0 \\ \text{or} \\ \text{if } h_\theta(x) < 0.5, y=1 \end{array} \\ 0 & \text{otherwise} \end{cases} \} \text{error}$$

$$\text{test error} = \frac{1}{m_{\text{test}}} \sum_{i=1}^{m_{\text{test}}} \text{err} (h_\theta(x_{\text{test}}^{(i)}), y^{(i)})$$

this is exactly the fraction of the examples  
in my test set that my hypothesis has mislabeled

# Summary

Once an hypothesis is done and the parameters that minimise the training error have been extracted:

- firstly, trouble shoots for errors in your predictions by getting more training examples, trying smaller sets of features, etc..
- secondly, move on to evaluating your hypothesis

Still, a hypothesis might have a low error for the training examples but still be inaccurate - because of overfitting.

Suggestion (simplest!) is to split the dataset into a **training set** and a **test set** - with a typical 70/30 ratio.

The new procedure is:

- Learn  $\Theta$  and minimise  $J_{\text{train}}(\Theta)$  - that we called just  $J(\Theta)$  in the previous slides
- compute the test set error  $J_{\text{test}}(\Theta)$ 
  - ❖ for linear regression: see formula
  - ❖ for classification: see the 0/1 misclassification error (proportion of the test data that has been badly misclassified)

The **test error** for the test set gives you a score, "how good I was", in predicting a value or in performing a classification

11100011101001001100100101001110101001001010111001  
0101001010101010110100101010101110001110100100110011  
01001010011101010010010101110010101001010111000110101  
**10110 Advices 00 for 01 applying 00 ML** 1010101011101010010101  
00110101001100100101011101010010101010101101001010  
101011110001110100100110010010100111010110110101110  
01010100101010101101000101010111000111010010011001  
10100101001110100100101011001010010011101011010110110  
101110010100101010101101000101010111000111001010  
01111101110 **Evaluating 10 a 10 learning 01 model 00 100100110101**  
110010100101010101101000101010111000111001000111  
**111011110 Model 10 selection 01 and 11 TVT 10 sets** 10110101001110  
101100110101110010100101010101011010001010101111000  
11100101000110111000100110011010010100111010101011  
100010010110010101010101110100011101010100111010111  
01101010101100001110000101100011010110101110

---

---

We have discussed the perils of overfitting, in which just because a learning algorithm fits a training set well, that does not mean it is a good hypothesis.

More generally, this is why the training set's error is not a good predictor for how well the hypothesis will do on new examples. And, similarly speaking,  $J_{\text{train}}(\Theta)$  is unlikely to be a good estimate of your actual generalisation error: it is likely to be lower than the actual generalisation error.

And, for this purpose, we introduced a **test** set.

Key question now: **is it enough?**

## model selection

$$d=1 \quad h_{\theta}(x) = \theta_0 + \theta_1 x$$

① take  $d=1$ , minimize training error  
⇒ you get some parameter vector  $\theta^{(1)}$

$$d=2 \quad h_{\theta}(x) = \theta_0 + \theta_1 x + \theta_2 x^2$$

② same with  $d=2 \Rightarrow \theta^{(2)}$

$$d=3 \quad h_{\theta}(x) = \theta_0 + \theta_1 x + \theta_2 x^2 + \theta_3 x^3 \quad ③ \Rightarrow \theta^{(3)}$$

...

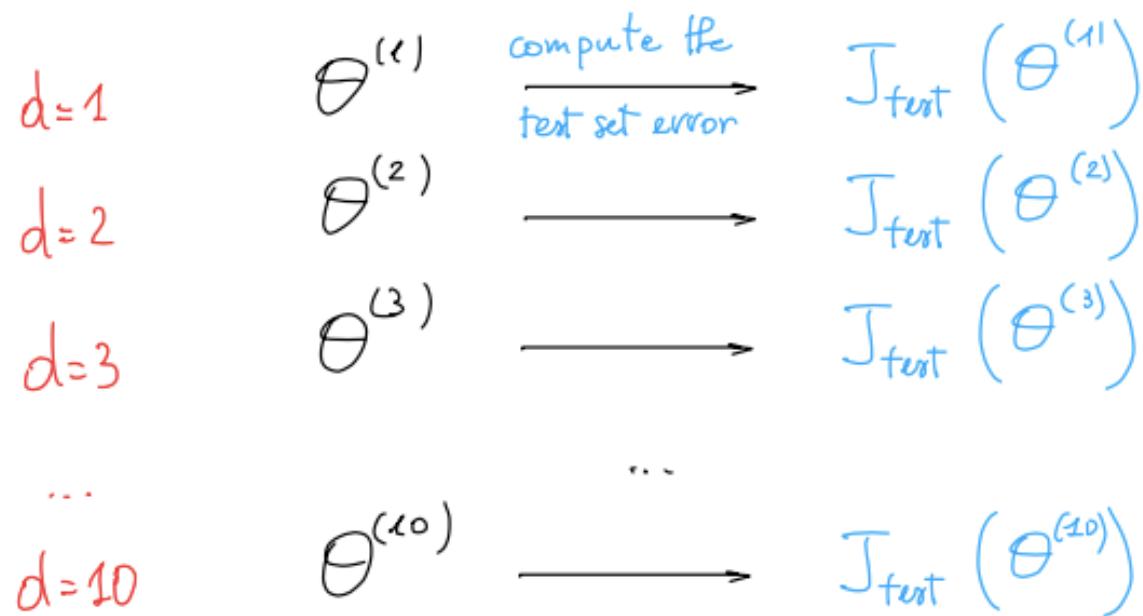
...

$$d=10 \quad h_{\theta}(x) = \theta_0 + \theta_1 x + \theta_2 x^2 + \dots + \theta_{10} x^{10}$$

⑩  $\Rightarrow \theta^{(10)}$

Notation :  $\theta^{(i)}$  = parameters vector I get by fitting polynomial of degree  $i$  to the training data

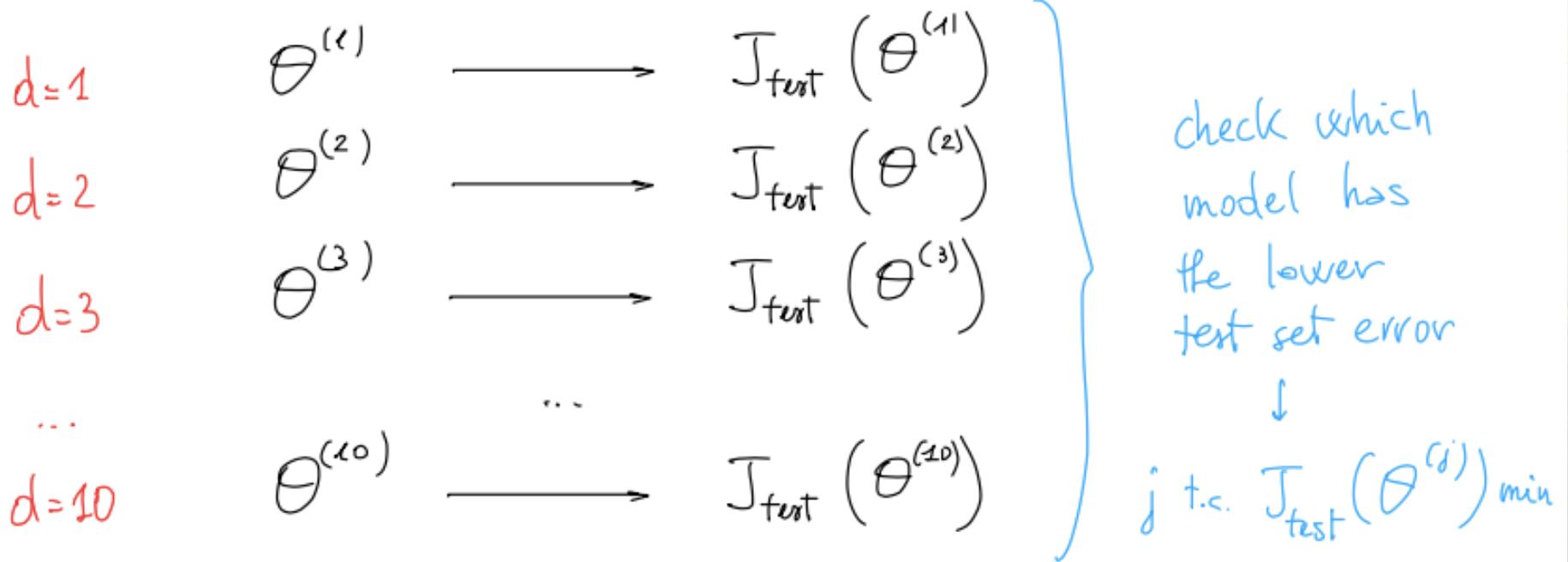
## model selection



Notation :  $\theta^{(i)}$  = parameters vector I get by fitting polynomial of degree  $i$  to the training data

$J_{\text{test}}(\theta^{(i)})$  = test set error - a measure of the performance on the test set

## model selection



## model selection

$$\begin{array}{lll} d=1 & \theta^{(1)} & \xrightarrow{\quad} J_{\text{test}}(\theta^{(1)}) \\ d=2 & \theta^{(2)} & \xrightarrow{\quad} J_{\text{test}}(\theta^{(2)}) \\ d=3 & \theta^{(3)} & \xrightarrow{\quad} J_{\text{test}}(\theta^{(3)}) \\ \dots & \dots & \dots \\ d=10 & \theta^{(10)} & \xrightarrow{\quad} J_{\text{test}}(\theta^{(10)}) \end{array}$$

Suppose  $j=5$  is  $j$  t.c.  $J_{\text{test}}(\theta^{(j)})$  minimal

$$\Rightarrow h_{\theta}(x) = \theta_0 + \theta_1 x + \theta_2 x^2 + \theta_3 x^3 + \theta_4 x^4 + \theta_5 x^5$$

How well does this model generalize?

## model selection

$$\begin{array}{lll} d=1 & \theta^{(1)} & \longrightarrow J_{\text{test}}(\theta^{(1)}) \\ d=2 & \theta^{(2)} & \longrightarrow J_{\text{test}}(\theta^{(2)}) \\ d=3 & \theta^{(3)} & \longrightarrow J_{\text{test}}(\theta^{(3)}) \\ \dots & & \dots \\ d=10 & \theta^{(10)} & \longrightarrow J_{\text{test}}(\theta^{(10)}) \end{array}$$

How well does this model generalize?

Look at how well my 5<sup>th</sup> order polynomial hypothesis has done on my test set, i.e.  $J_{\text{test}}(\theta^{(5)})$

This will not be a fair estimate of how well my hypothesis generalizes : we chose  $d=5$  because it gave best  $J_{\text{test}}$  ; it is likely to be an optimistic estimate of the generalization error (namely, I have chosen  $d=5$  based on the test set !) so my hypothesis is likely to do better on the test set...

size (m <sup>2</sup> )	price (k€)
150	300
100	200
80	200
75	250
160	400
210	410
90	200
95	280
105	350
120	370

60%

20%

TRAINING  
SET

(CROSS-)VALIDATION  
SET

TEST  
SET

$(x^{(1)}, y^{(1)})$

⋮

$(x^{(m)}, y^{(m)})$

$m = \#$  training examples

$(x_{cv}^{(1)}, y_{cv}^{(1)})$

⋮

$(x_{cv}^{(m_{cv})}, y_{cv}^{(m_{cv})})$

$m_{cv} = \#$  cv examples

$(x_{test}^{(1)}, y_{test}^{(1)})$

⋮

$(x_{test}^{(m_{test})}, y_{test}^{(m_{test})})$

$m_{test} = \#$  test examples

---

---

So, now that we've defined the **training set**, (cross) **validation set** and **test set**..

.. we can also define the **training error**, (cross) **validation error**, and **test error**.

### TRAINING ERROR

$$J_{\text{train}}(\theta) = \frac{1}{2m_{\text{train}}} \sum_{i=1}^{m_{\text{train}}} \left( h_{\theta}(x_{\text{train}}^{(i)}) - y_{\text{train}}^{(i)} \right)^2$$

value as my previous  $J(\theta)$

### CROSS-VALIDATION ERROR

$$J_{\text{cv}}(\theta) = \frac{1}{2m_{\text{cv}}} \sum_{i=1}^{m_{\text{cv}}} \left( h_{\theta}(x_{\text{cv}}^{(i)}) - y_{\text{cv}}^{(i)} \right)^2$$

### TEST ERROR

$$J_{\text{test}}(\theta) = \frac{1}{2m_{\text{test}}} \sum_{i=1}^{m_{\text{test}}} \left( h_{\theta}(x_{\text{test}}^{(i)}) - y_{\text{test}}^{(i)} \right)^2$$

Similar formulas, but note there is **absolutely no overlap** among the  $m_{\text{train}}$ ,  $m_{\text{cv}}$  and  $m_{\text{test}}$  examples

# The procedure with training/validation/test (TVT)

model selection (using training /test /cv sets)

1.  $h_{\theta}(x) = \theta_0 + \theta_1 x \longrightarrow \min_{\theta} J(\theta) \rightarrow \theta^{(1)} \longrightarrow J_{cv}(\theta^{(1)})$

2.  $h_{\theta}(x) = \theta_0 + \theta_1 x + \theta_2 x^2 \longrightarrow \theta^{(2)} \longrightarrow J_{cv}(\theta^{(2)})$

3.  $h_{\theta}(x) = \theta_0 + \theta_1 x + \theta_2 x^2 + \theta_3 x^3 \longrightarrow \theta^{(3)} \longrightarrow J_{cv}(\theta^{(3)})$

...

10.  $h_{\theta}(x) = \theta_0 + \theta_1 x + \theta_2 x^2 + \dots + \theta_{10} x^{10} \rightarrow \theta^{(10)} \longrightarrow J_{cv}(\theta^{(10)})$

I test h on  
the cv set

Basically, we do not use the test set to select the model, but the CV set.

→  $J_{cv}(\theta^{(i)})$  tells me how well each h do on my CV set,

→ I pick the h with the lowest CV error

## model selection (using training / test / cv sets)

1.  $h_{\theta}(x) = \theta_0 + \theta_1 x \longrightarrow \theta^{(1)} \rightarrow J_{cv}(\theta^{(1)})$
2.  $h_{\theta}(x) = \theta_0 + \theta_1 x + \theta_2 x^2 \longrightarrow \theta^{(2)} \rightarrow J_{cv}(\theta^{(2)})$
3.  $h_{\theta}(x) = \theta_0 + \theta_1 x + \theta_2 x^2 + \theta_3 x^3 \longrightarrow \theta^{(3)} \rightarrow J_{cv}(\theta^{(3)})$
- ...
10.  $h_{\theta}(x) = \theta_0 + \theta_1 x + \theta_2 x^2 + \dots + \theta_{10} x^{10} \rightarrow \theta^{(10)} \rightarrow J_{cv}(\theta^{(10)})$

Suppose lowest is  $J_{cv}(\theta^{(4)})$ :

$$\Rightarrow h_{\theta}(x) = \theta_0 + \theta_1 x + \theta_2 x^2 + \theta_3 x^3 + \theta_4 x^4 \quad \text{good! selected!}$$

Now:

Use the test set to estimate the generalization error of the model  
 selected using the (CV set.)  $\rightarrow$  i.e. compute  $J_{test}(\theta^{(4)})$

## Take-away message

---

Take your dataset, split into TVT, isolate the **test set** (→ “Put it in the closet, hide the key!”). Then, perform your ML modelling with the rest of the datasets.

**NEVER, EVER LOOK AT THE TEST SET  
in this early phase.**

Bring back the key, get the test set back on the table, only when you are done with model creation, training, tuning.. and you use it ONLY to assess the generalisation error.

# Summary

A learning algo that fits a training set well is not necessarily a good hypothesis. It could overfit and, as a result, you would be entering the real application of your ML model being unmotivatedly optimistic, and your final predictions would eventually be disappointingly poor.

- the error of your  $h$  as measured on the training set will be lower than the error on any other data set

If you e.g. need to decide among many models, each with a different polynomial degree, in order to identify the “best” one you can compute the error result for each.

**Break your dataset into e.g. Train 60%, CV 20%, Test 20%**

Then, follow this procedure:

- optimise the  $\theta$ s using the training set (polynomial  $h$ ? then, try each polynomial degree  $d$ )
- find the best model (polynomial? find the best  $d$ ), i.e. smallest error, using the cv set
- estimate the generalisation error using the test set (polynomial? with  $J_{\text{test}}(\Theta^{(d)})$ )

In this way (polynomial example),  $d$  has not been trained using the test set.

# A final remark (personal opinion)

In the code, all this will be very simple and quick, and often almost automatic to execute in your ML projects implementation, and you will be greatly helped by libraries to do all this.

But:

That does not mean this part is unnecessary. Quite the opposite. It gives solidity to the steps you would do quickly afterwards. It is important, IMO, in a basic applied ML course, that you spend some time to think about the **diligence needed in dealing with your dataset and the TVT splitting**: it will increase your awareness on what you do, your ability to spot logical mistakes, your speed in debugging code, your talent in spotting data leakage evidences, etc.

11100011101001001100100101001110101001001010111001  
0101001010101010110100101010101110001110100100110011  
01001010011101010010010101110010101001010111000110101  
**10110 Advices 00 for 01 applying 00 ML** 101010101110101010010101  
0011010100110010010101110101010010101010101101001010  
1010111100011101001001100100101001110101101101101110  
01010100101010101101000101010111000111010010011001  
10100101001110100100101011001010100100111010110110110  
101110010100101010101101000101010111000111001010  
01111101110 **Bias 10 and 10 Variance** 0100100100101101110101  
1110010100101010101101000101010111100011100101001  
**1111101110 Diagnosis 01 bias 10 vs 10 variance** 01001001001001110  
10110011010111001010100101010101011010001010101111000  
1110010100011011110001001100110100101001110101011101011  
10001001011001010101010111010001110101010011101011101011  
0110101010101100001110000101100011010110101110

---

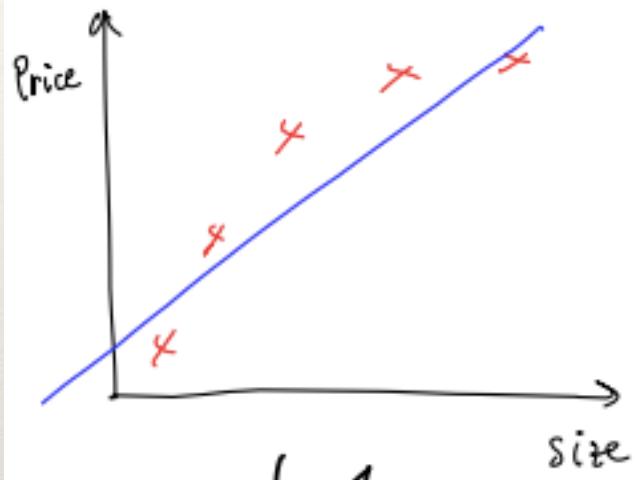
---

If we run a learning algorithm to make predictions and it does not do as well as we hoped, almost always it is because we have either a **high bias problem** or a **high variance problem**

- in other words, either an **underfitting** problem or an **overfitting** problem

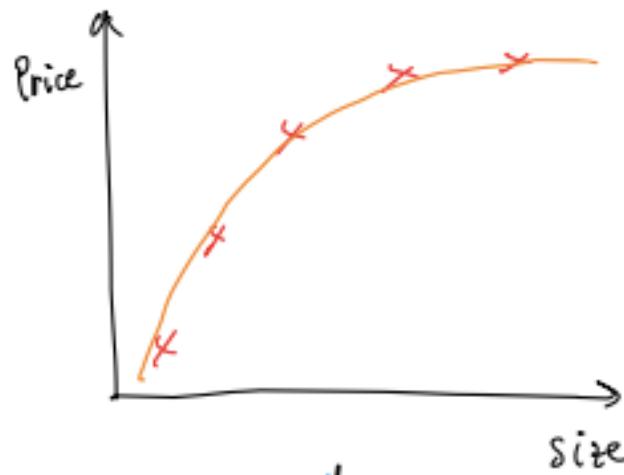
Let's dive more into this.

Restart from this :



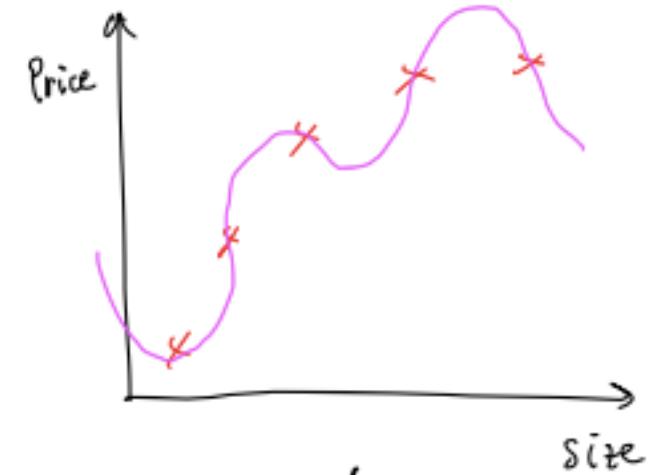
$$\theta_0 + \theta_1 x$$

"UNDERFITTING"  
or  
"HIGH-BIAS"



$$\theta_0 + \theta_1 x + \theta_2 x^2$$

"JUST RIGHT"



$$\theta_0 + \theta_1 x + \theta_2 x^2 + \theta_3 x^3 + \theta_4 x^4$$

"OVERFITTING"  
or  
"HIGH-VARIANCE"

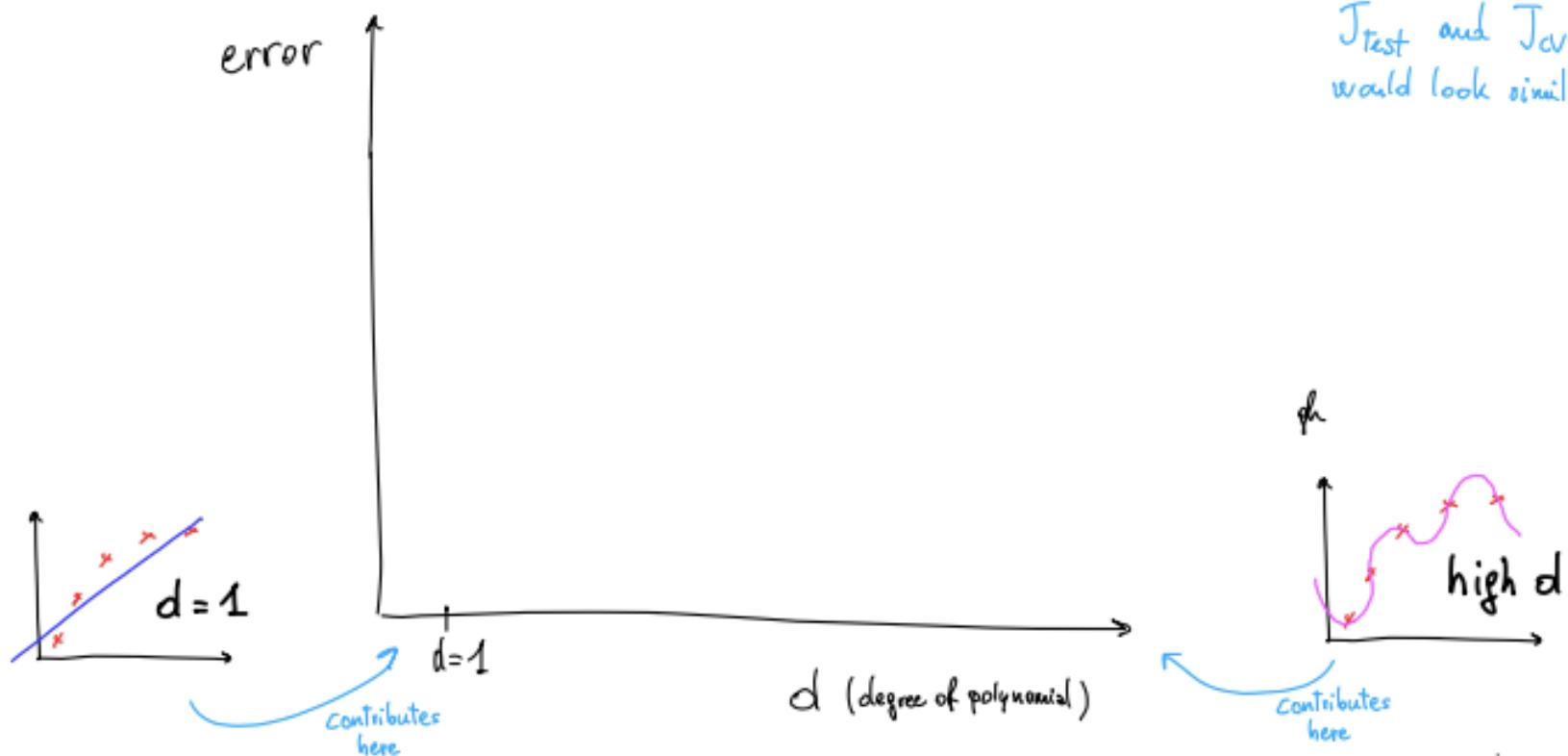
### TRAINING ERROR

$$J_{\text{train}}(\theta) = \frac{1}{2m_{\text{train}}} \sum_{i=1}^m \left( h_{\theta}(x_{\text{train}}^{(i)}) - y_{\text{train}}^{(i)} \right)^2$$

### CROSS-VALIDATION ERROR

$$J_{\text{cv}}(\theta) = \frac{1}{2m_{\text{cv}}} \sum_{i=1}^{m_{\text{cv}}} \left( h_{\theta}(x_{\text{cv}}^{(i)}) - y_{\text{cv}}^{(i)} \right)^2$$

NOTE:  
 $J_{\text{test}}$  and  $J_{\text{cv}}$   
 would look similar

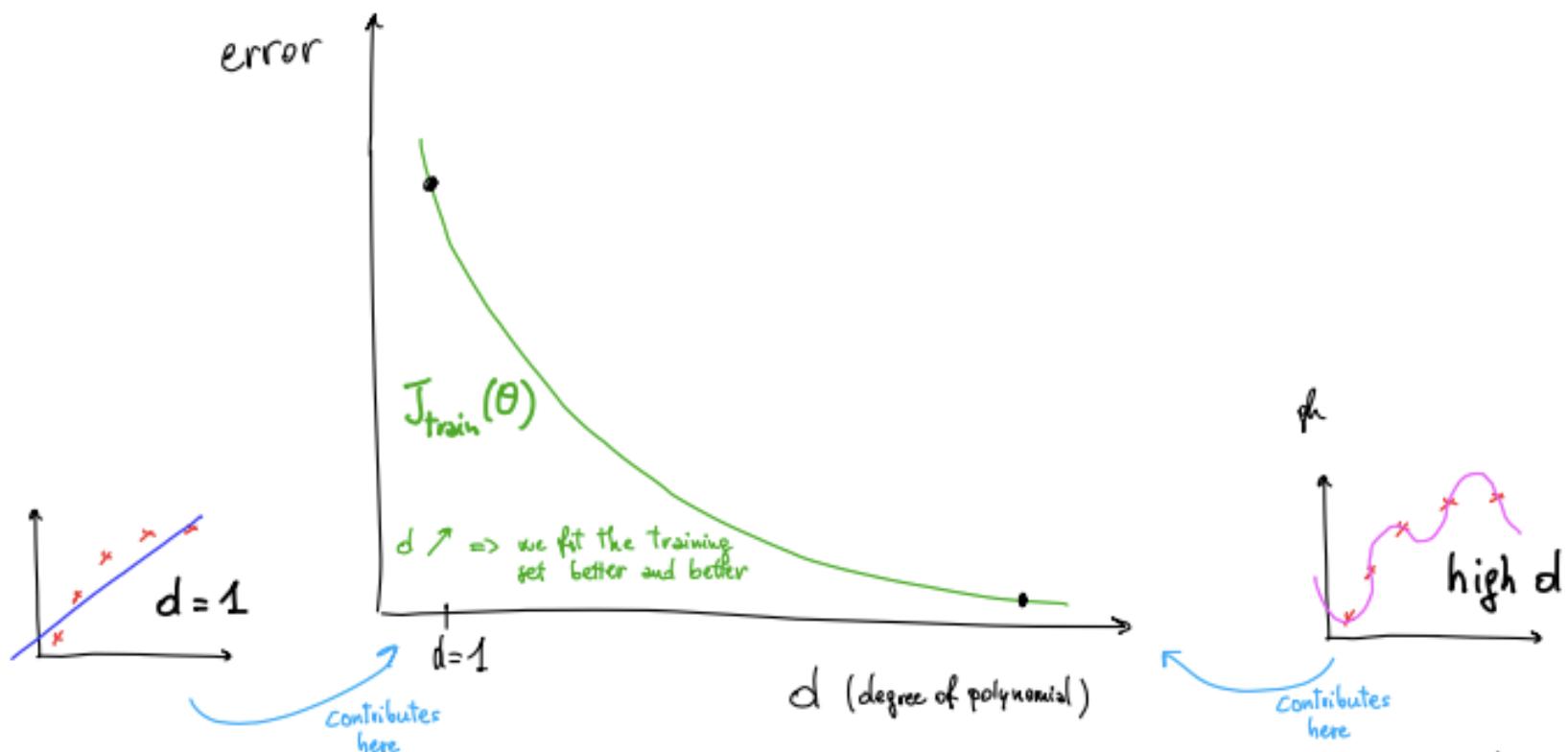


### TRAINING ERROR

$$J_{\text{train}}(\theta) = \frac{1}{2m_{\text{train}}} \sum_{i=1}^m \left( h_{\theta}(x_{\text{train}}^{(i)}) - y_{\text{train}}^{(i)} \right)^2$$

### CROSS-VALIDATION ERROR

$$J_{\text{cv}}(\theta) = \frac{1}{2m_{\text{cv}}} \sum_{i=1}^{m_{\text{cv}}} \left( h_{\theta}(x_{\text{cv}}^{(i)}) - y_{\text{cv}}^{(i)} \right)^2$$



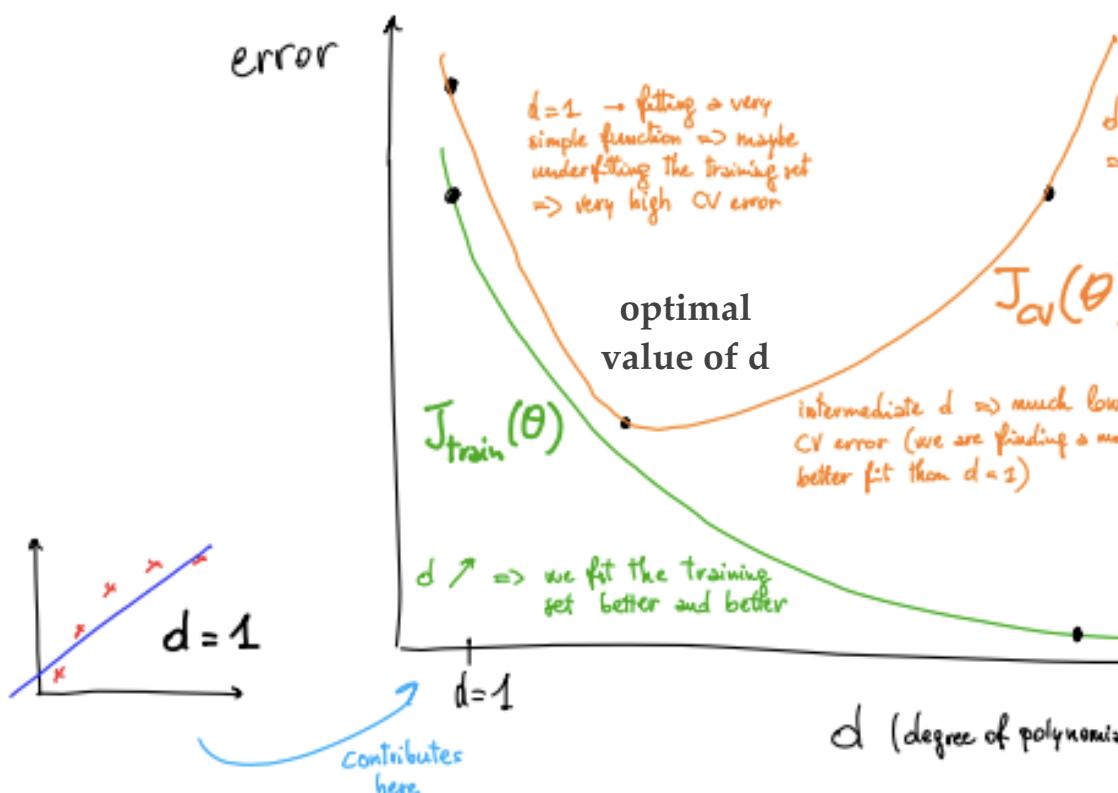
### TRAINING ERROR

$$J_{\text{train}}(\theta) = \frac{1}{2m_{\text{train}}} \sum_{i=1}^{m_{\text{train}}} (h_{\theta}(x_{\text{train}}^{(i)}) - y_{\text{train}}^{(i)})^2$$

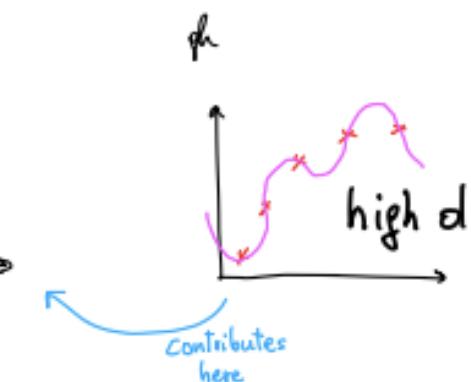
### CROSS-VALIDATION ERROR

$$J_{\text{cv}}(\theta) = \frac{1}{2m_{\text{cv}}} \sum_{i=1}^{m_{\text{cv}}} (h_{\theta}(x_{\text{cv}}^{(i)}) - y_{\text{cv}}^{(i)})^2$$

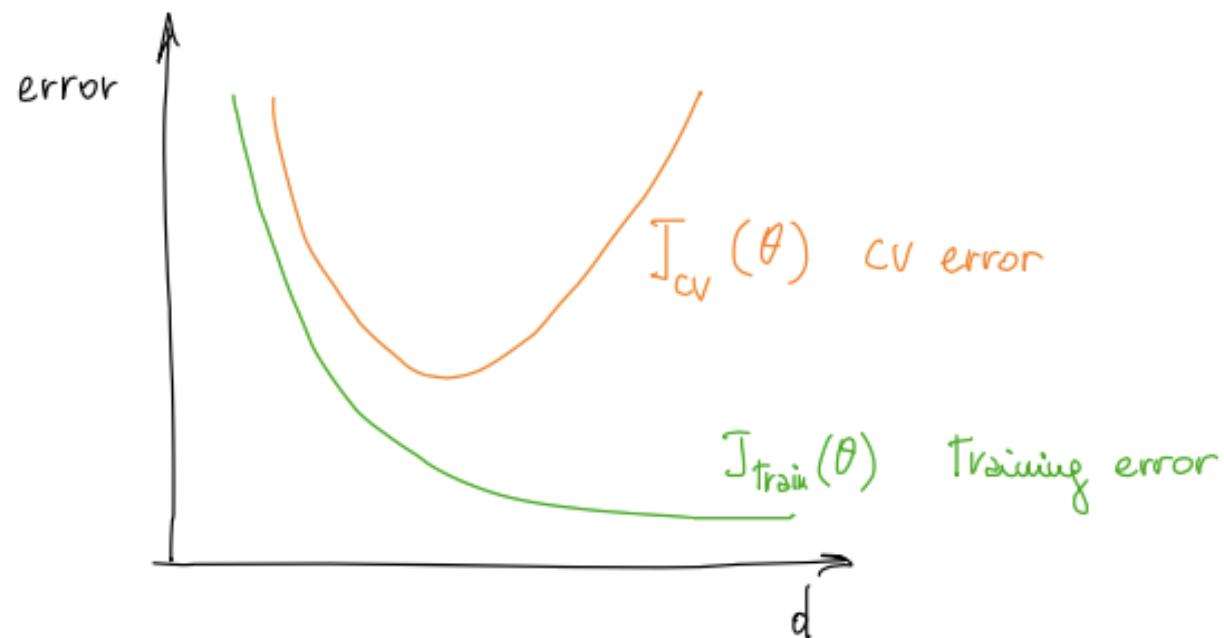
NOTE:  
 $J_{\text{test}}$  and  $J_{\text{cv}}$   
 would look similar



$J_{\text{cv}}$  (or  $J_{\text{test}}$ ) do not  
 continue to decrease



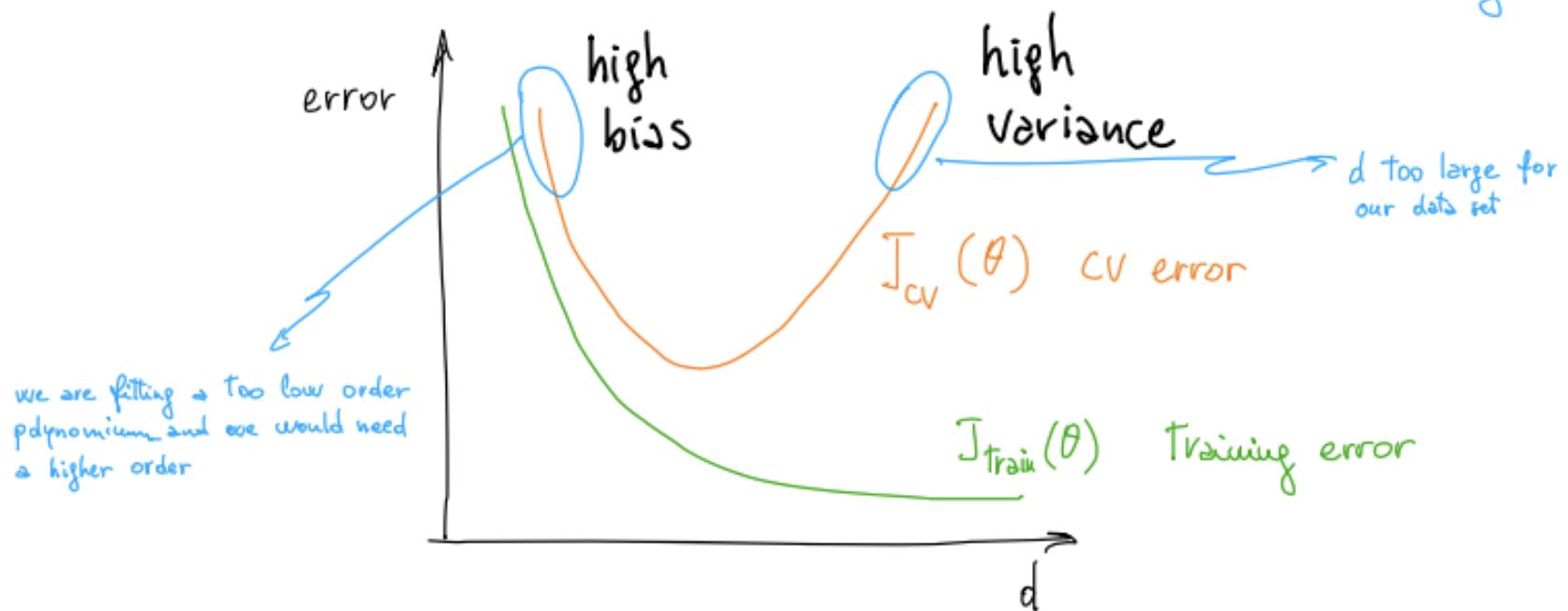
Suppose your learning algo is performing less well than hoped.  
Is it a bias problem or a variance problem?



Suppose your learning algo is performing less well than hoped.

Is it a bias problem or a variance problem?  
(high bias) (high variance)

$J_{cv}(\theta)$  (or  $J_{test}(\theta)$ )  
too high



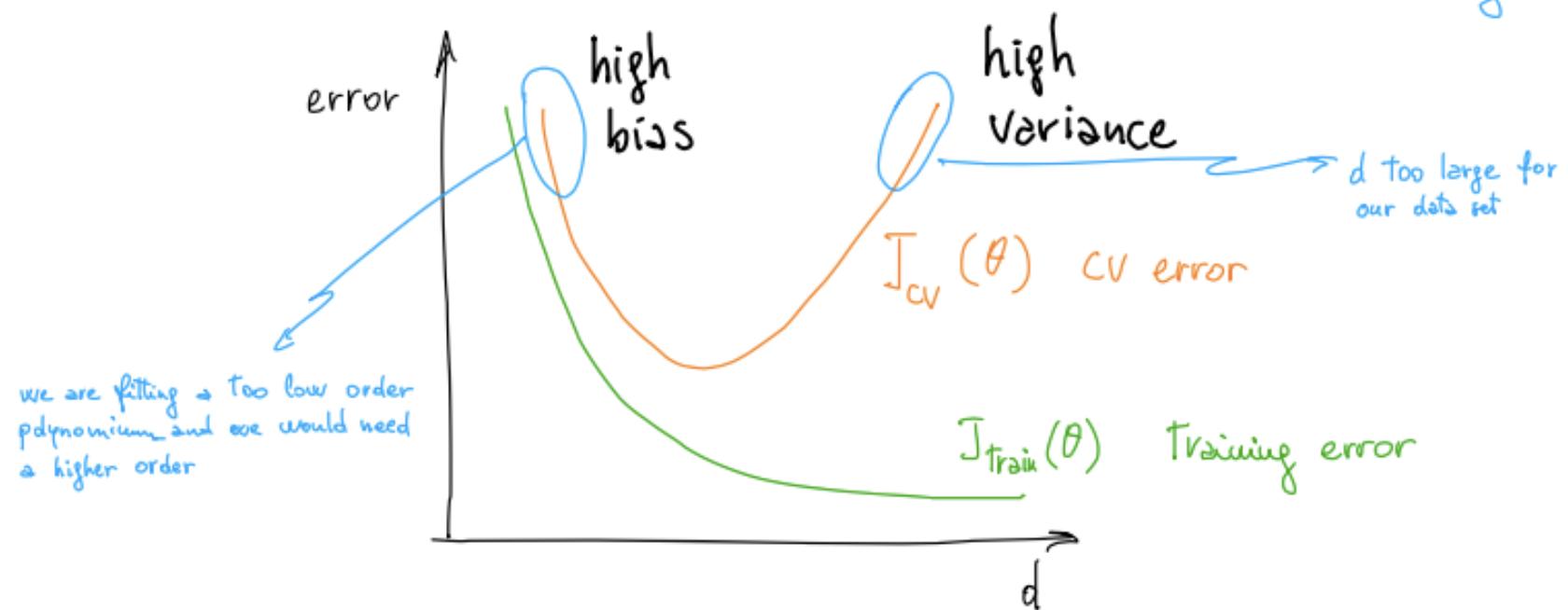
Suppose your learning algo is performing less well than hoped.

Is it a bias problem or a variance problem?

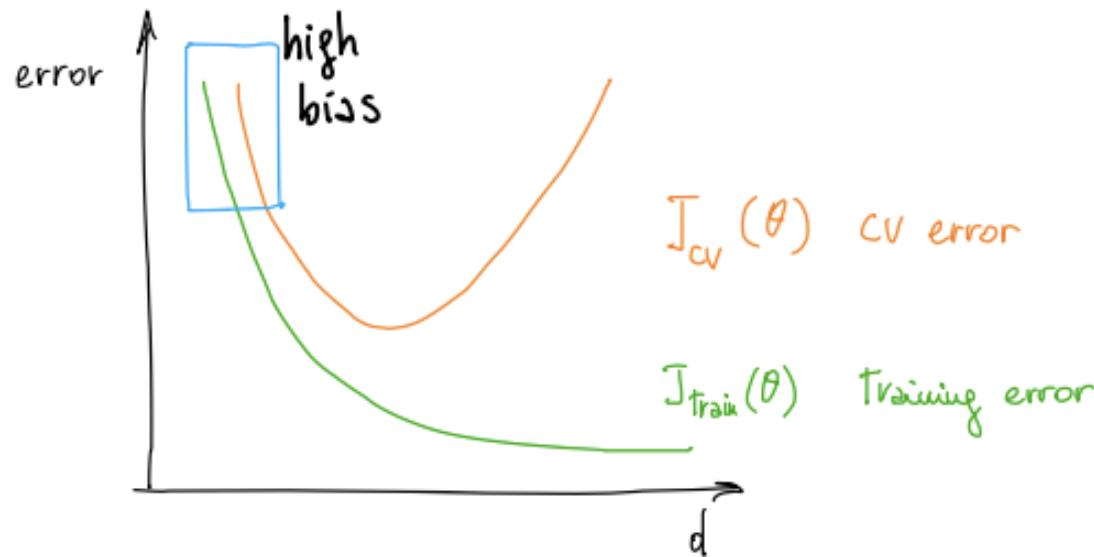
(high bias)

(high variance)

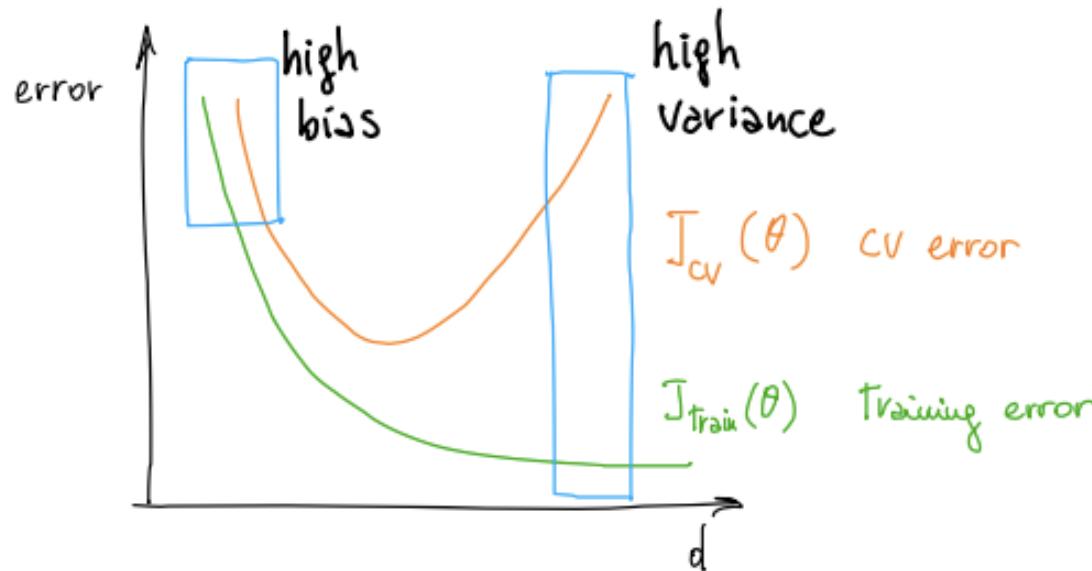
$J_{cv}(\theta)$  (or  $J_{test}(\theta)$ )  
too high



This gives us a way to disentangle high bias vs high variance!



BIAS (underfit)  $\rightarrow$   $J_{\text{train}}(\theta)$  high (h not fitting training set well)  
 $J_{\text{cv}}(\theta)$  high ( $\approx J_{\text{train}}$ , maybe higher)



BIAS (underfit)  $\rightarrow$   $J_{\text{train}}(\theta)$  high ( $h$  not fitting training set well)  
 $J_{\text{cv}}(\theta)$  high ( $\approx J_{\text{train}}$ , maybe higher)

VARIANCE (overfit)  $\rightarrow$   $J_{\text{train}}(\theta)$  low (ie. fitting training set well)  
 $J_{\text{cv}}(\theta)$  high ( $\gg J_{\text{train}}$ )

# Summary

We discussed the **diagnosis of bias vs variance** in case of bad predictions

We examined, in our explanatory example, the relationship between the degree of the polynomial  $d$  and the underfitting vs overfitting of our hypothesis.

We associate **high bias** to underfitting, and **high variance** to overfitting. We need to distinguish whether bias or variance is the problem contributing to bad predictions. Ideally, we need to find a golden working point between these 2 points

The training error will tend to decrease as we increase the degree  $d$  of the polynomial. At the same time, the cv error will tend to decrease as we increase  $d$  up to a point, and then it will increase as  $d$  is increasing, forming a convex curve

**High bias (underfitting):**

- both  $J_{\text{train}}(\Theta)$  and  $J_{\text{cv}}(\Theta)$  will be high. Also,  $J_{\text{cv}}(\Theta) \approx J_{\text{train}}(\Theta)$

**High variance (overfitting):**

- $J_{\text{train}}(\Theta)$  will be low and  $J_{\text{cv}}(\Theta)$  will be  $>> J_{\text{train}}(\Theta)$

11100011101001001100100101001110101001001010111001  
0101001010101010110100101010101110001110100100110011  
01001010011101010010010101110010101001010111000110101  
**10110 Advices 00 for 01 applying 00 ML** 101010101110101010010101  
0011010100110010010101110101010010101010101101001010  
1010111100011101001001100100101001110101101101101110  
01010100101010101101000101010111000111010010011001  
1010010100111010010010110010101001001110101101101110  
101110010100101010101101000101010111000111001010  
01111101110 **Bias 10 and 10 Variance** 0100100100101101110101  
111001010010101010110100010101011100011100101001  
1111101110 **Regularisation 01 and 01 bias 1/0 variance** 010100111010  
110011010111001010100101010101010110100010101011100011  
10010101000110111100010011001001010011101010101110  
001001011001010101010111010001110101010011101011101  
10101010101100001110000101100011010110101110

---

---

We have seen how regularisation can help prevent overfitting.

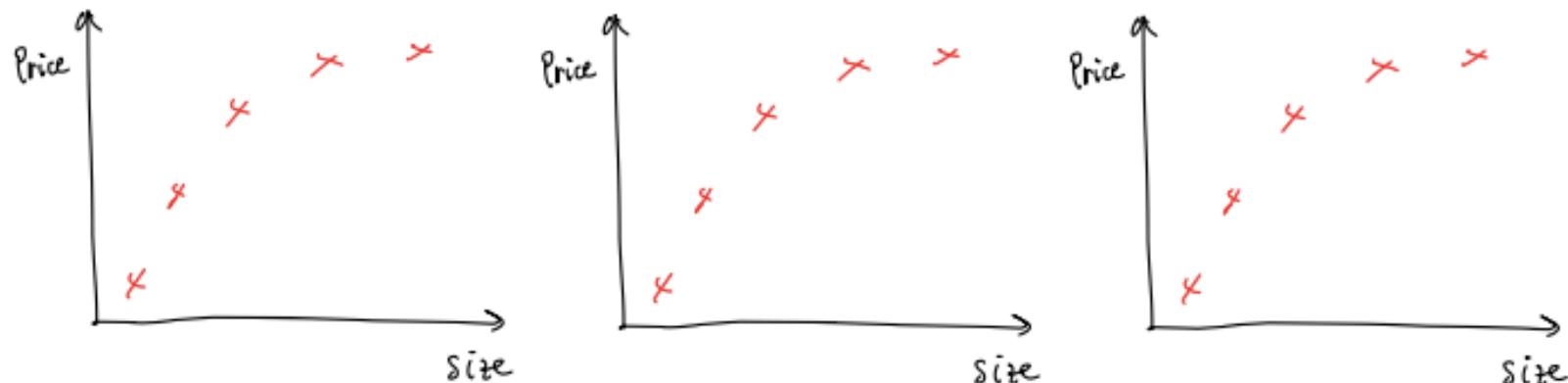
How does it affect (visibly?) the bias vs variances of a learning algorithm, that we discussed?

Linear regression (w regularization)

already applied to attack overfitting

$$h_{\theta}(x) = \theta_0 + \theta_1 x + \theta_2 x^2 + \theta_3 x^3 + \theta_4 x^4$$

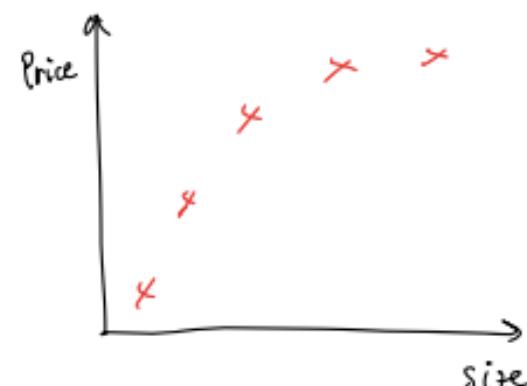
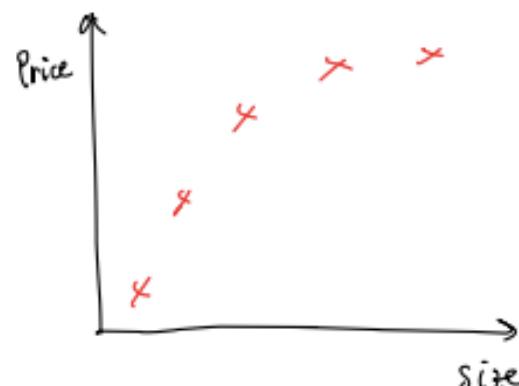
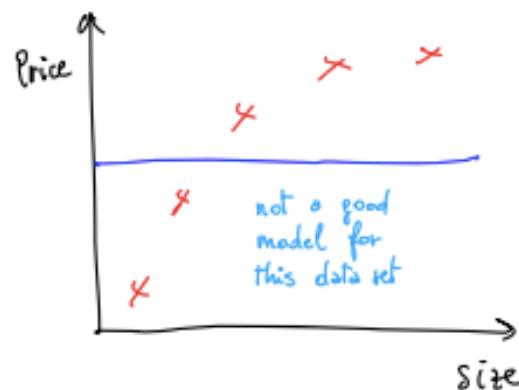
$$J(\theta) = \frac{1}{2m} \sum_{i=1}^m (h_{\theta}(x^{(i)}) - y^{(i)})^2 + \frac{\lambda}{2m} \sum_{j=1}^m \theta_j^2$$



Linear regression (w regularization) already applied to attack overfitting

$$h_{\theta}(x) = \theta_0 + \theta_1 x + \theta_2 x^2 + \theta_3 x^3 + \theta_4 x^4$$

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



large  $\lambda$  penalized

$$\lambda = 10'000 \quad \theta_1, \dots, \theta_n \approx 0$$

$$h_{\theta}(x) \approx \theta_0$$

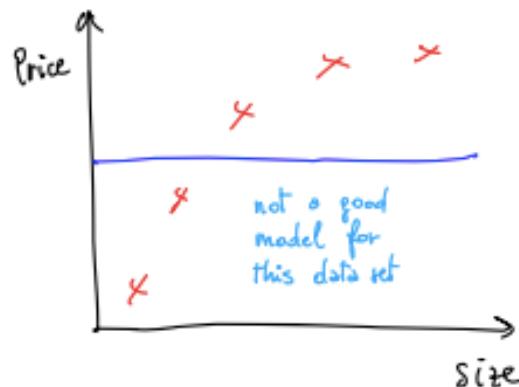
high bias (underfitting)

## Linear regression (w regularization)

already applied to attack overfitting

$$h_{\theta}(x) = \theta_0 + \theta_1 x + \theta_2 x^2 + \theta_3 x^3 + \theta_4 x^4$$

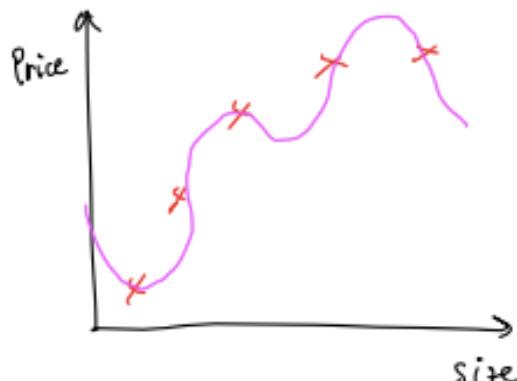
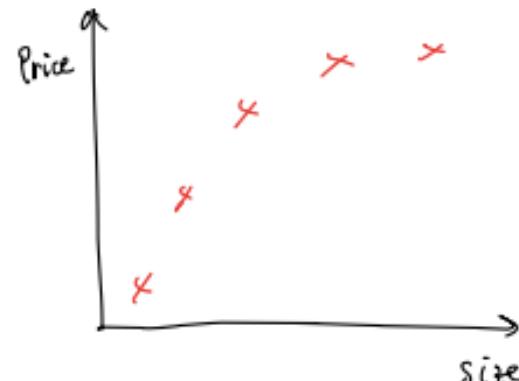
$$J(\theta) = \frac{1}{2m} \sum_{i=1}^m (h_{\theta}(x^{(i)}) - y^{(i)})^2 + \frac{\lambda}{2m} \sum_{j=1}^m \theta_j^2$$



large  $\lambda$  penalized

$$\lambda = 10'000 \quad \theta_1, \dots, \theta_n \approx 0 \\ h_{\theta}(x) \approx \theta_0$$

high bias (underfitting)



small  $\lambda$

$$\lambda = 0 \quad h_{\theta}(x) \text{ high d}$$

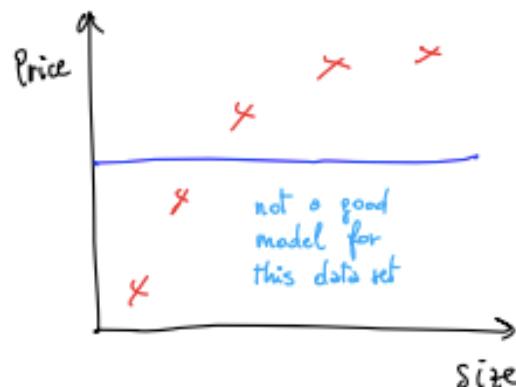
high variance (overfitting)  
fitting w/o regularization

## Linear regression (w regularization)

already applied to attack overfitting

$$h_{\theta}(x) = \theta_0 + \theta_1 x + \theta_2 x^2 + \theta_3 x^3 + \theta_4 x^4$$

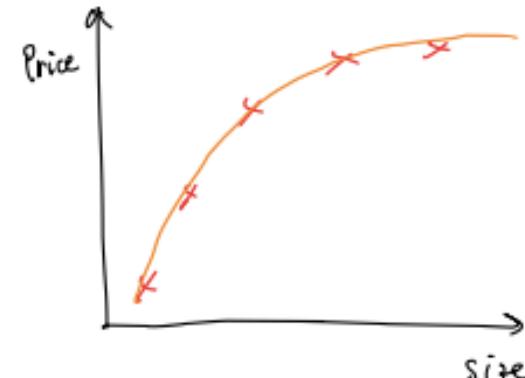
$$J(\theta) = \frac{1}{2m} \sum_{i=1}^m (h_{\theta}(x^{(i)}) - y^{(i)})^2 + \frac{\lambda}{2m} \sum_{j=1}^m \theta_j^2$$



large  $\lambda$

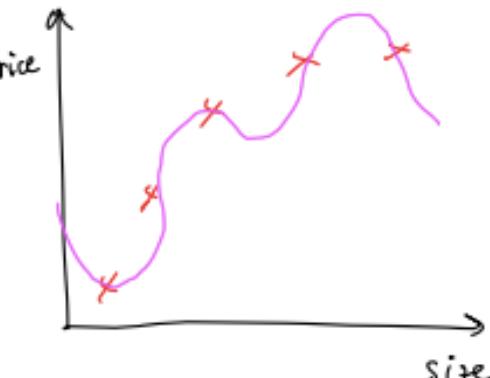
$$\lambda = 10'000 \quad \theta_1, \dots, \theta_n \approx 0 \quad h_{\theta}(x) \approx \theta_0$$

high bias (underfitting)



intermediate  $\lambda$

"just right"



small  $\lambda$

$$\lambda = 0 \quad h_{\theta}(x) \text{ high d}$$

high variance (overfitting)  
fitting w/o regularization

---

---

How do we automatically choose a good value for the regularisation parameter  $\lambda$ ?

We want to apply the model selection procedure (seen before) to selecting the regularisation parameter  $\lambda$ .

Choosing the regularization parameter  $\lambda$  (automatically!)

$$h_{\theta}(x) = \theta_0 + \theta_1 x + \theta_2 x^2 + \theta_3 x^3 + \theta_4 x^4$$

$$J(\theta) = \frac{1}{2m} \sum_{i=1}^n (h_{\theta}(x^{(i)}) - y^{(i)})^2 + \frac{\lambda}{2m} \sum_{j=1}^m \theta_j^2$$

Choosing the regularization parameter  $\lambda$  (automatically!)

$$h_{\theta}(x) = \theta_0 + \theta_1 x + \theta_2 x^2 + \theta_3 x^3 + \theta_4 x^4$$

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

Let's define  $J_{\text{train}}$ ,  $J_{\text{cv}}$ ,  $J_{\text{test}}$  to be the optimization objective but w/o the regular. term (AS BEFORE)

$$J_{\text{train}}(\theta) = \frac{1}{2m_{\text{train}}} \sum_{i=1}^m (h_{\theta}(x_{\text{train}}^{(i)}) - y_{\text{train}}^{(i)})^2$$

$$J_{\text{cv}}(\theta) = \frac{1}{2m_{\text{cv}}} \sum_{i=1}^{m_{\text{cv}}} (h_{\theta}(x_{\text{cv}}^{(i)}) - y_{\text{cv}}^{(i)})^2$$

$$J_{\text{test}}(\theta) = \frac{1}{2m_{\text{test}}} \sum_{i=1}^{m_{\text{test}}} (h_{\theta}(x_{\text{test}}^{(i)}) - y_{\text{test}}^{(i)})^2$$

## Choosing the regularization parameter $\lambda$ (automatically!)

$$h_{\theta}(x) = \theta_0 + \theta_1 x + \theta_2 x^2 + \theta_3 x^3 + \theta_4 x^4$$

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

1. try  $\lambda = 0$  ← not using regal.

2. try  $\lambda = 0.01$

3. try  $\lambda = 0.02$

4. try  $\lambda = 0.04$

5. try  $\lambda = 0.08$

...

12. try  $\lambda = 10$   
 $(\sim 10.24)$

e.g.  
 multiplier of 2

e.g.  
 12 different models,  
 each with a different  $\lambda$

Choosing the regularization parameter  $\lambda$  (automatically!)

$$h_{\theta}(x) = \theta_0 + \theta_1 x + \theta_2 x^2 + \theta_3 x^3 + \theta_4 x^4$$

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

1. try  $\lambda = 0 \rightarrow \min_{\theta} J(\theta) \rightarrow \theta^{(1)}$

2. try  $\lambda = 0.01 \rightarrow \min_{\theta} J(\theta) \rightarrow \theta^{(2)}$

3. try  $\lambda = 0.02 \dashrightarrow \theta^{(3)}$

4. try  $\lambda = 0.04 \dots$

5. try  $\lambda = 0.08$

...

12. try  $\lambda = 10 \dashrightarrow \theta^{(12)}$

## Choosing the regularization parameter $\lambda$ (automatically!)

$$h_{\theta}(x) = \theta_0 + \theta_1 x + \theta_2 x^2 + \theta_3 x^3 + \theta_4 x^4$$

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

1. try  $\lambda = 0$

$$\min_{\theta} J(\theta) \rightarrow \theta^{(1)}$$

2. try  $\lambda = 0.01$

$$\min_{\theta} J(\theta) \rightarrow \theta^{(2)}$$

3. try  $\lambda = 0.02$

$$\dots \rightarrow \theta^{(3)}$$

4. try  $\lambda = 0.04$

...

5. try  $\lambda = 0.08$

...

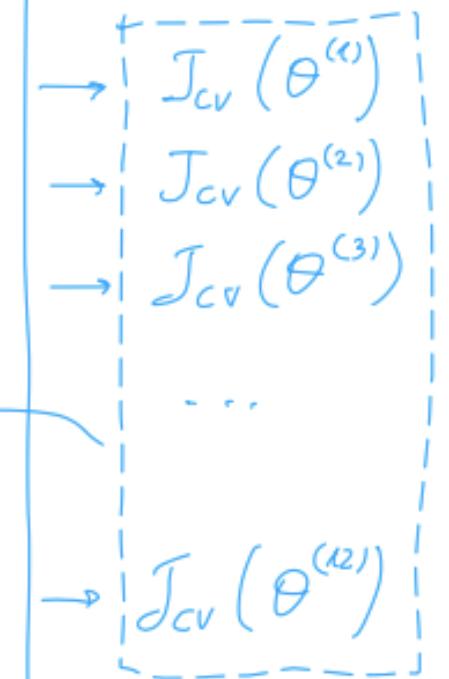
12. try  $\lambda = 10$

$$\rightarrow \theta^{(12)}$$

pick model  $j$  that gives  
smallest  $J_{cv}(\theta^{(j)}) \Rightarrow$   
you get your  $\lambda$

first step

evaluate them with CV set



second step

Choosing the regularization parameter  $\lambda$  (automatically!)

$$h_{\theta}(x) = \theta_0 + \theta_1 x + \theta_2 x^2 + \theta_3 x^3 + \theta_4 x^4$$

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

1. try  $\lambda = 0$

2. try  $\lambda = 0.01$

3. try  $\lambda = 0.02$

4. try  $\lambda = 0.04$

5. try  $\lambda = 0.08$

...

12. try  $\lambda = 10$



(...)  $\Rightarrow$



lowest is e.g.  $J_{cv}(\theta^{(5)})$

Report the test set error  
for the  $\theta^{(5)}$  choice, i.e.:

$J_{test}(\theta^{(5)})$

---

---

So, this is model selection method applied also to **selection of the best regularization parameter  $\lambda$**

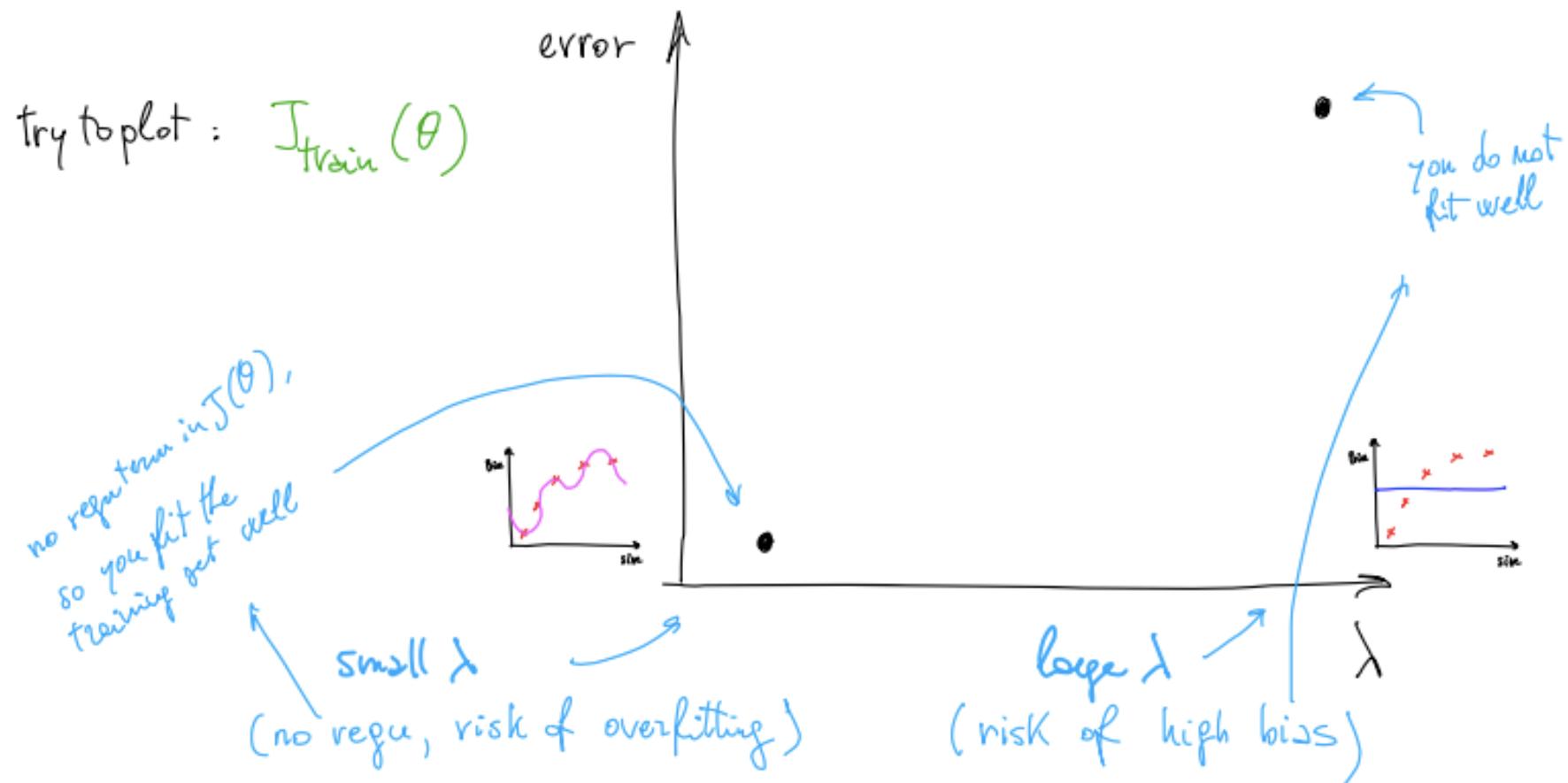
- it is the logic behind the implementation of various hyper parameters scan method you will just import and use in your code!

One last thing we should cover is to get a better understanding of how cross validation and training errors vary as we vary the regularisation parameter  $\lambda$ .

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

$$\rightarrow J_{\text{train}}(\theta) = \frac{1}{2m_{\text{train}}} \sum_{i=1}^{m_{\text{train}}} (h_\theta(x_{\text{train}}^{(i)}) - y_{\text{train}}^{(i)})^2$$

$$\rightarrow J_{\text{cv}}(\theta) = \frac{1}{2m_{\text{cv}}} \sum_{i=1}^{m_{\text{cv}}} (h_\theta(x_{\text{cv}}^{(i)}) - y_{\text{cv}}^{(i)})^2$$



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

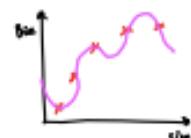
$$\rightarrow J_{\text{train}}(\theta) = \frac{1}{2m_{\text{train}}} \sum_{i=1}^m (h_\theta(x_{\text{train}}^{(i)}) - y_{\text{train}}^{(i)})^2$$

$$\rightarrow J_{\text{cv}}(\theta) = \frac{1}{2m_{\text{cv}}} \sum_{i=1}^{m_{\text{cv}}} (h_\theta(x_{\text{cv}}^{(i)}) - y_{\text{cv}}^{(i)})^2$$

try to plot:  $J_{\text{train}}(\theta)$

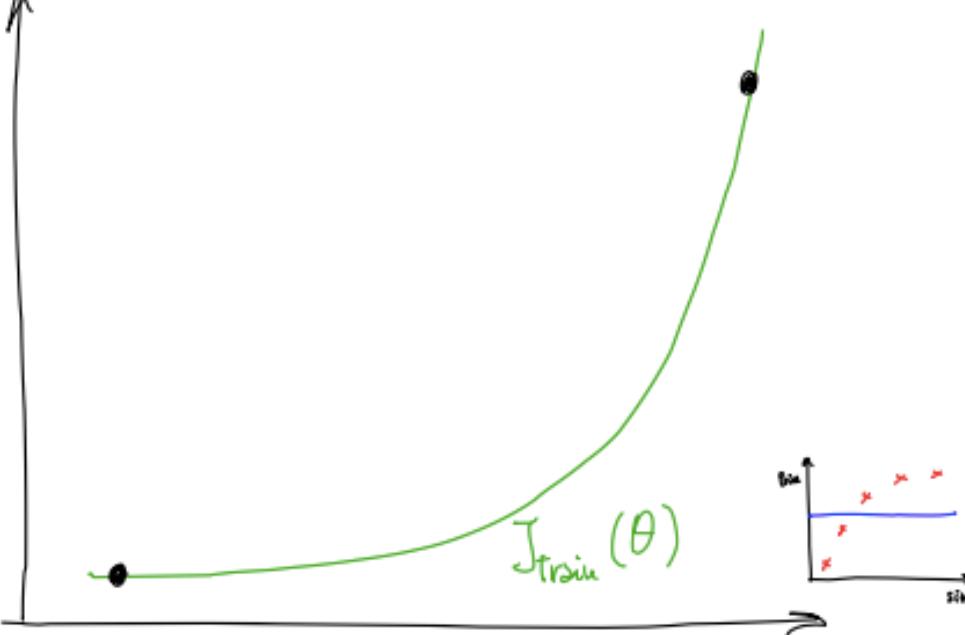
~~DONE!~~

error



small  $\lambda$

(no regu, risk of overfitting)



large  $\lambda$

(risk of high bias)

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

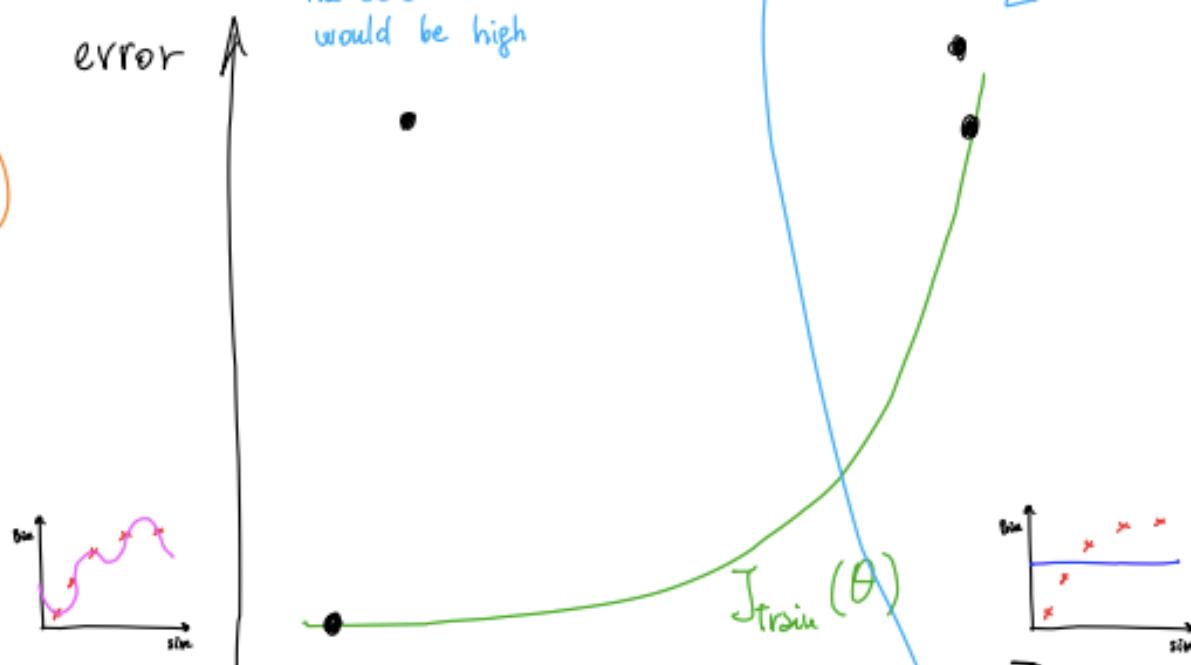
$$\rightarrow J_{\text{train}}(\theta) = \frac{1}{2m_{\text{train}}} \sum_{i=1}^{m_{\text{train}}} (h_\theta(x^{(i)}_{\text{train}}) - y^{(i)}_{\text{train}})^2$$

$$\rightarrow J_{\text{CV}}(\theta) = \frac{1}{2m_{\text{cv}}} \sum_{i=1}^{m_{\text{cv}}} (h_\theta(x^{(i)}_{\text{cv}}) - y^{(i)}_{\text{cv}})^2$$

high variance regime,  
we may be overfitting  
the data  $\Rightarrow$  CV error  
would be high

we underfit  $\Rightarrow$  bias regime  
 $\Rightarrow$  CV error would be high,  
and we would not do well  
on the CV test

try to plot:  $J_{\text{CV}}(\theta)$



small  $\lambda$

(no regu, risk of overfitting)

large  $\lambda$

(risk of high bias)

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

$$\rightarrow J_{\text{train}}(\theta) = \frac{1}{2m_{\text{train}}} \sum_{i=1}^m (h_\theta(x_{\text{train}}^{(i)}) - y_{\text{train}}^{(i)})^2$$

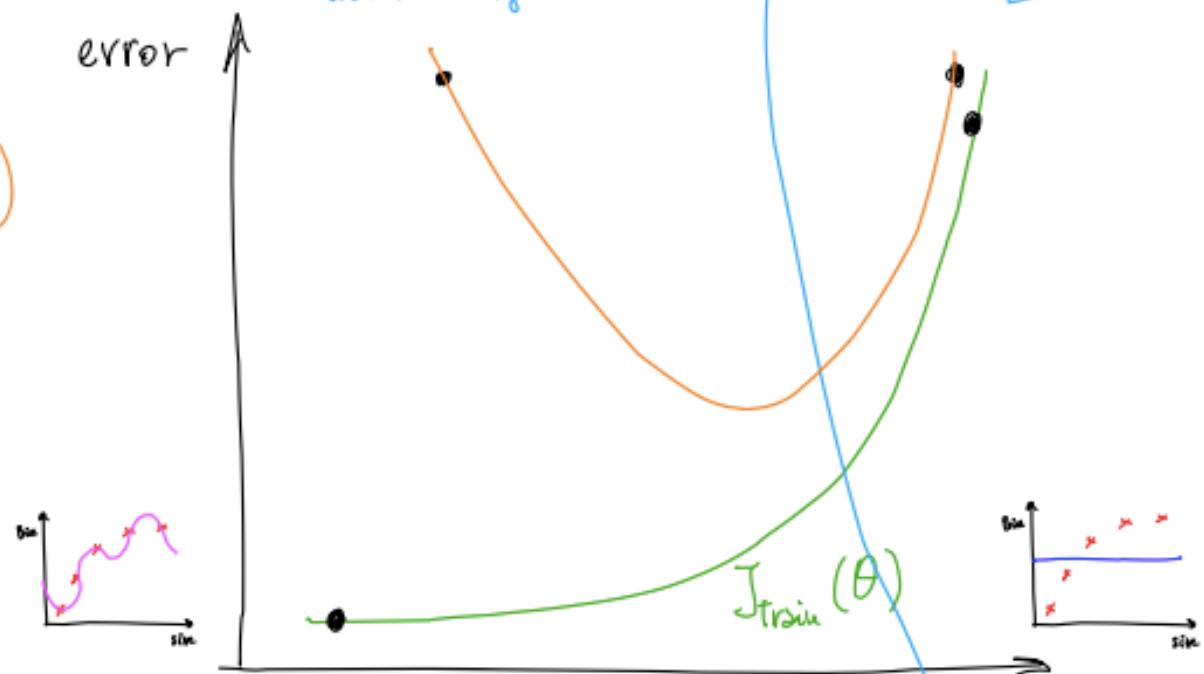
$$\rightarrow J_{\text{cv}}(\theta) = \frac{1}{2m_{\text{cv}}} \sum_{i=1}^{m_{\text{cv}}} (h_\theta(x_{\text{cv}}^{(i)}) - y_{\text{cv}}^{(i)})^2$$

high variance regime,  
we may be overfitting  
the data  $\Rightarrow$  CV error  
would be high

we underfit  $\Rightarrow$  bias regime  
 $\Rightarrow$  CV error would be high,  
and we would not do well  
on the CV test

try to plot:  $J_{\text{cv}}(\theta)$

DONE!

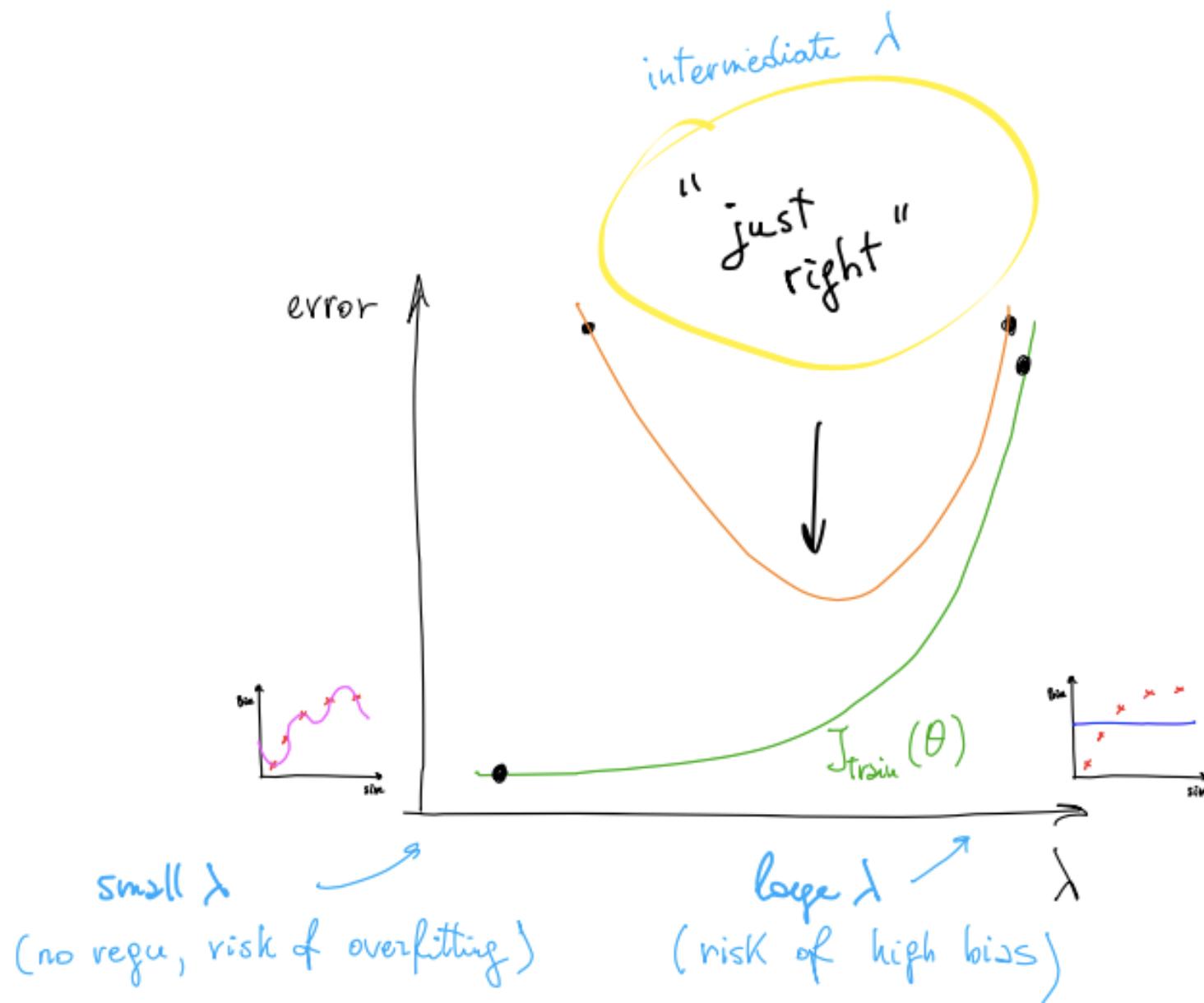


small  $\lambda$

(no regu, risk of overfitting)

large  $\lambda$

(risk of high bias)



# Summary (1/2)

If we plot our “price to pay” versus  $\lambda$  we see that:

- as  $\lambda$  increases, our fit becomes more rigid
- on the other hand, as  $\lambda$  approaches 0, we tend to overfit the data

So, how do we choose the best  $\lambda$  to be “just right”?

In order to choose the model and the regularisation term  $\lambda$  we could, as a principle, to:

- create a list of  $\lambda$ s (i.e. 0, 0.01, 0.02, ..., 10)
- create a set of models with different degrees of any other variants
- iterate through the  $\lambda$ s and for each  $\lambda$  go through all the models and learn  $\Theta$ s
- compute the cv error using the learned  $\Theta$  (computed with a  $\lambda$ ) on the  $J_{cv}(\Theta)$  without regularisation (i.e.  $\lambda=0$ )
- select the best combination  $(\Theta, \lambda)$  that produces the lowest error on the cv set
- use  $(\Theta, \lambda)$ , apply it to  $J_{test}(\Theta)$  to see if it has a good generalisation of the problem

## Summary (2/2)

---

Note: the drawn curves are somewhat cartoon-like; real curves are not so idealised, and more messy and noisy

But:

Trends are meaningful. Sometimes on a dataset you spot these trends if you plot a cv error, and then you can (either manually or automatically) select a point that minimises the cv error and yields a value of  $\lambda$  **that** corresponds to a low cv error - i.e. presumably also a low generalisation error.

In other words, plots like these help to understand better what's going on, and help to verify that you are picking a good value for the regularisation parameter  $\lambda$ .

---

# Next

---

We are ready to put together a diagnostic tool called **learning curves**.

11100011101001001100100101001110101001001010111001  
0101001010101010110100101010101110001110100100110011  
01001010011101010010010101110010101001010111000110101  
**10110 Advices 00 for 01 applying 00 ML** 101010101110101010010101  
0011010100110010010101110101010010101010101101001010  
1010111100011101001001100100101001110101101101101110  
01010100101010101101000101010111000111010010011001  
10100101001110100100101011001010100100111010110110110  
101110010100101010101101000101010111000111001010  
01111101110 **Bias 10 and 10 Variance** 0100100100101101110101  
111001010010101010110100010101011100011100101001  
**1111101110 Learning 01 curves** 01011101011011101010010100111  
0101100110101110010101001010101010110100010101011100  
011100101000110111100010011001010010100111010101  
110001001011001010101010111010001110101010011101011  
10110101010101100001110000101100011010110101110

---

---

If either you wanted a sanity check that your algorithm is working correctly, or if you want to improve the performance of the algorithm, **learning curves** are often very useful to plot.

It is a tool often used to diagnose if a learning algorithm may be suffering from a bias problem or a variance problem, or a little bit of both.

# Learning curves

$$J_{\text{train}}(\theta) = \frac{1}{2m_{\text{train}}} \sum_{i=1}^m (h_{\theta}(x_{\text{train}}^{(i)}) - y_{\text{train}}^{(i)})^2$$
$$J_{\text{cv}}(\theta) = \frac{1}{2m_{\text{cv}}} \sum_{i=1}^{m_{\text{cv}}} (h_{\theta}(x_{\text{cv}}^{(i)}) - y_{\text{cv}}^{(i)})^2$$

Plot now as a function of  $m$  (training set size)

start with better ;

$$J_{\text{train}}(\theta) = \frac{1}{2m_{\text{train}}} \sum_{i=1}^m (h_{\theta}(x_{\text{train}}^{(i)}) - y_{\text{train}}^{(i)})^2$$

$$J_{\text{cv}}(\theta) = \frac{1}{2m_{\text{cv}}} \sum_{i=1}^{m_{\text{cv}}} (h_{\theta}(x_{\text{cv}}^{(i)}) - y_{\text{cv}}^{(i)})^2$$

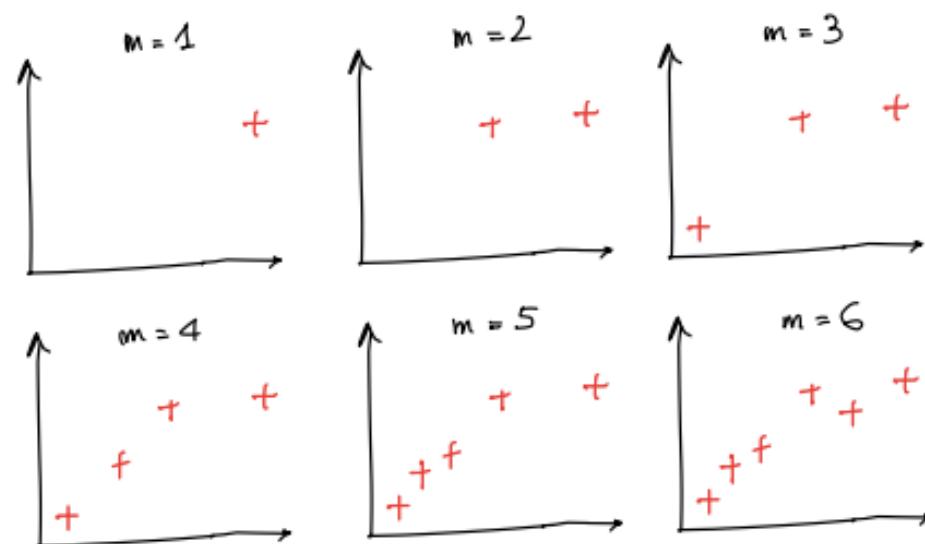
start with  
later

Plot now as a function of  $m$  (training set size)

To do so, start with  $m=1$  and add examples one by one ...

Suppose 6 examples, and you try this:

$$h_{\theta}(x) = \theta_0 + \theta_1 x + \theta_2 x^2$$



$$J_{\text{train}}(\theta) = \frac{1}{2m_{\text{train}}} \sum_{i=1}^m \left( h_{\theta}(x_{\text{train}}^{(i)}) - y_{\text{train}}^{(i)} \right)^2$$

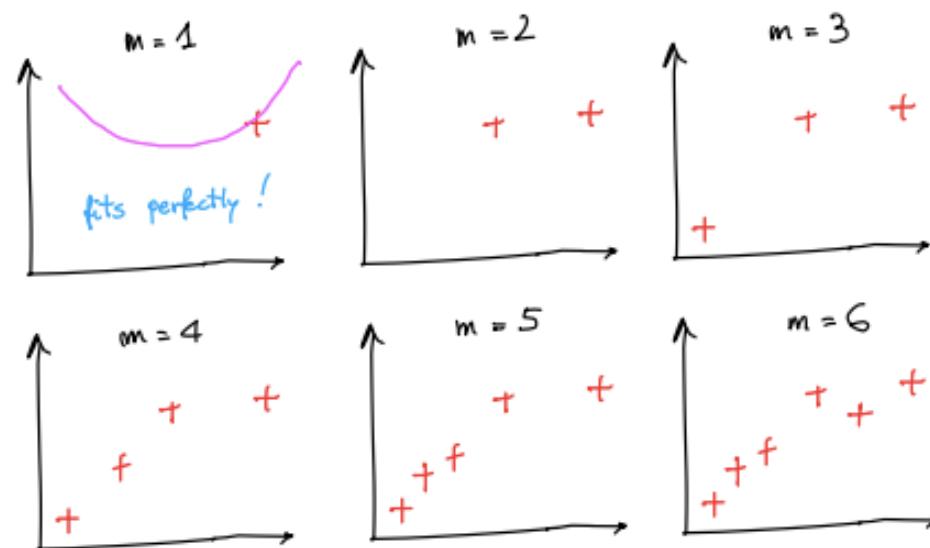
$$J_{\text{cv}}(\theta) = \frac{1}{2m_{\text{cv}}} \sum_{i=1}^{m_{\text{cv}}} \left( h_{\theta}(x_{\text{cv}}^{(i)}) - y_{\text{cv}}^{(i)} \right)^2$$

Plot now as a function of  $m$  (training set size)

To do so, start with  $m=1$  and add examples one by one ...

Suppose 6 examples, and you try this:

$$h_{\theta}(x) = \theta_0 + \theta_1 x + \theta_2 x^2$$



$$J_{\text{train}}(\theta) = \frac{1}{2m_{\text{train}}} \sum_{i=1}^m \left( h_{\theta}(x_{\text{train}}^{(i)}) - y_{\text{train}}^{(i)} \right)^2$$

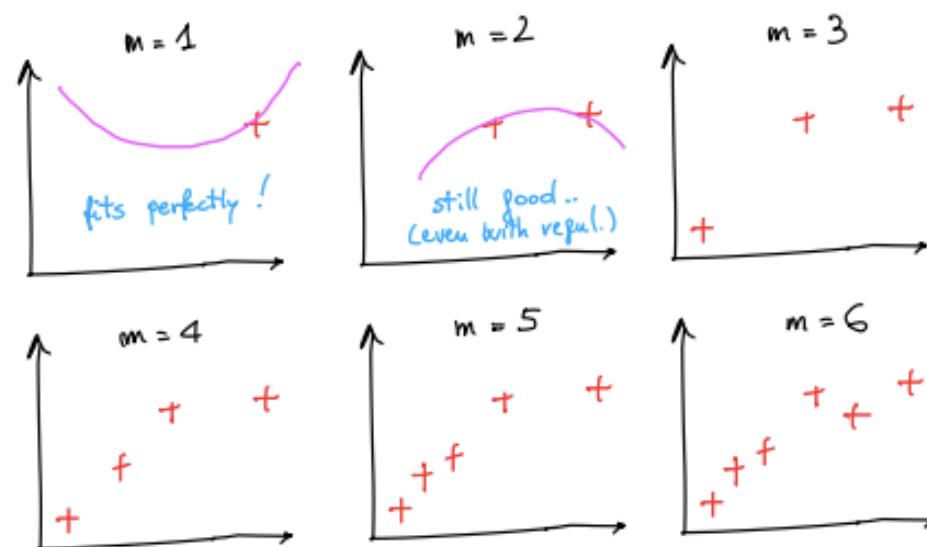
$$J_{\text{cv}}(\theta) = \frac{1}{2m_{\text{cv}}} \sum_{i=1}^{m_{\text{cv}}} \left( h_{\theta}(x_{\text{cv}}^{(i)}) - y_{\text{cv}}^{(i)} \right)^2$$

Plot now as a function of  $m$  (training set size)

To do so, start with  $m=1$  and add examples one by one ...

Suppose 6 examples, and you try this:

$$h_{\theta}(x) = \theta_0 + \theta_1 x + \theta_2 x^2$$



$$J_{\text{train}}(\theta) = \frac{1}{2m_{\text{train}}} \sum_{i=1}^m \left( h_{\theta}(x_{\text{train}}^{(i)}) - y_{\text{train}}^{(i)} \right)^2$$

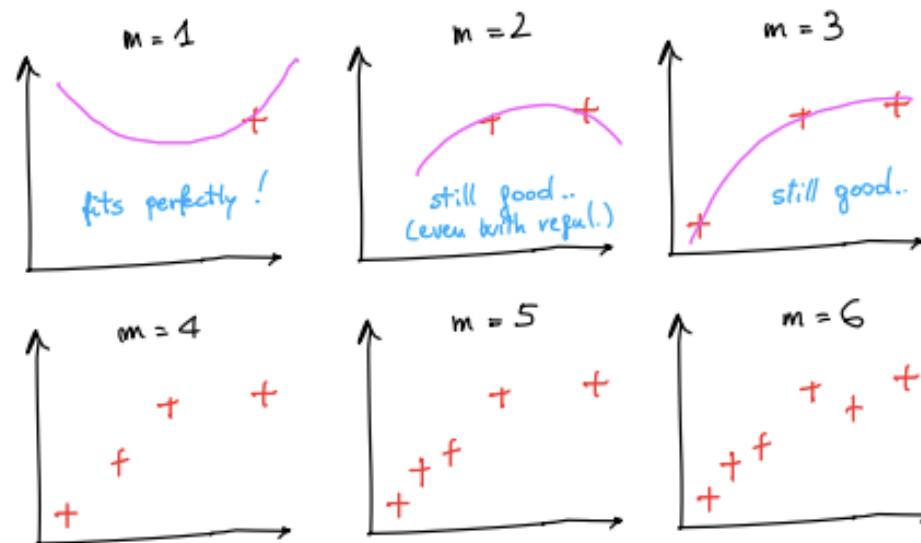
$$J_{\text{cv}}(\theta) = \frac{1}{2m_{\text{cv}}} \sum_{i=1}^{m_{\text{cv}}} \left( h_{\theta}(x_{\text{cv}}^{(i)}) - y_{\text{cv}}^{(i)} \right)^2$$

Plot now as a function of  $m$  (training set size)

To do so, start with  $m=1$  and add examples one by one ...

Suppose 6 examples, and you try this:

$$h_{\theta}(x) = \theta_0 + \theta_1 x + \theta_2 x^2$$



$$J_{\text{train}}(\theta) = \frac{1}{2m_{\text{train}}} \sum_{i=1}^m \left( h_{\theta}(x_{\text{train}}^{(i)}) - y_{\text{train}}^{(i)} \right)^2$$

$$J_{\text{cv}}(\theta) = \frac{1}{2m_{\text{cv}}} \sum_{i=1}^{m_{\text{cv}}} \left( h_{\theta}(x_{\text{cv}}^{(i)}) - y_{\text{cv}}^{(i)} \right)^2$$

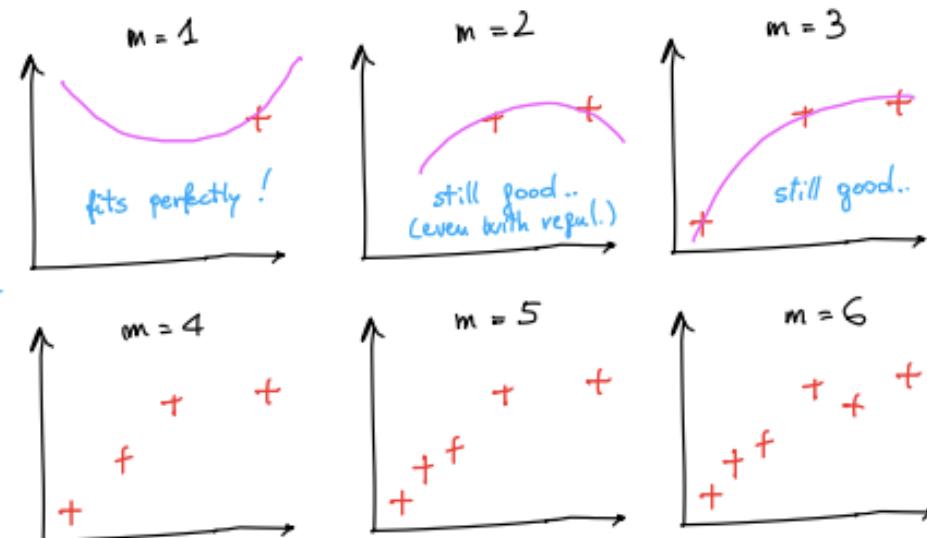
Plot now as a function of  $m$  (training set size)

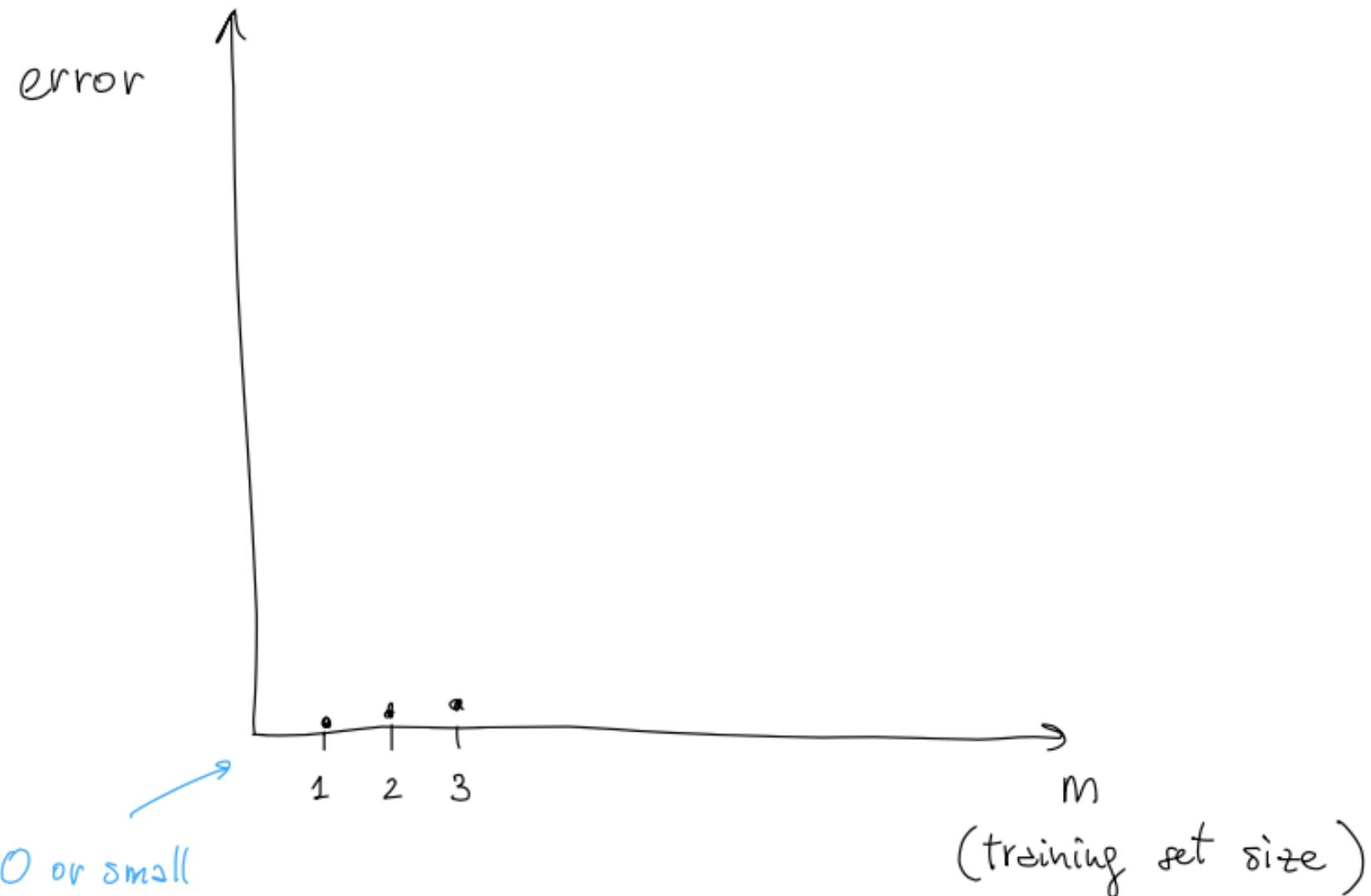
To do so, start with  $m=1$  and add examples one by one ...

Suppose 6 examples, and you try this:

$$h_{\theta}(x) = \theta_0 + \theta_1 x + \theta_2 x^2$$

$m=1, 2, 3 \Rightarrow$  training error w/o regul.  
is 0 (or slightly higher than - but close  
to - w regul.)





$$J_{\text{train}}(\theta) = \frac{1}{2m_{\text{train}}} \sum_{i=1}^m \left( h_{\theta}(x_{\text{train}}^{(i)}) - y_{\text{train}}^{(i)} \right)^2$$

$$J_{\text{cv}}(\theta) = \frac{1}{2m_{\text{cv}}} \sum_{i=1}^{m_{\text{cv}}} \left( h_{\theta}(x_{\text{cv}}^{(i)}) - y_{\text{cv}}^{(i)} \right)^2$$

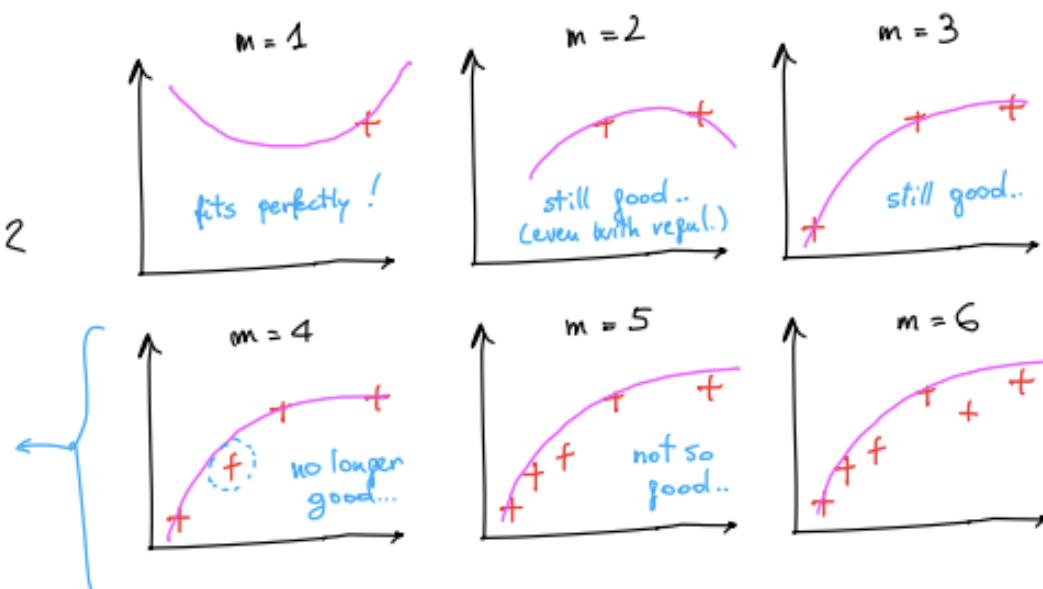
Plot now as a function of  $m$  (training set size)

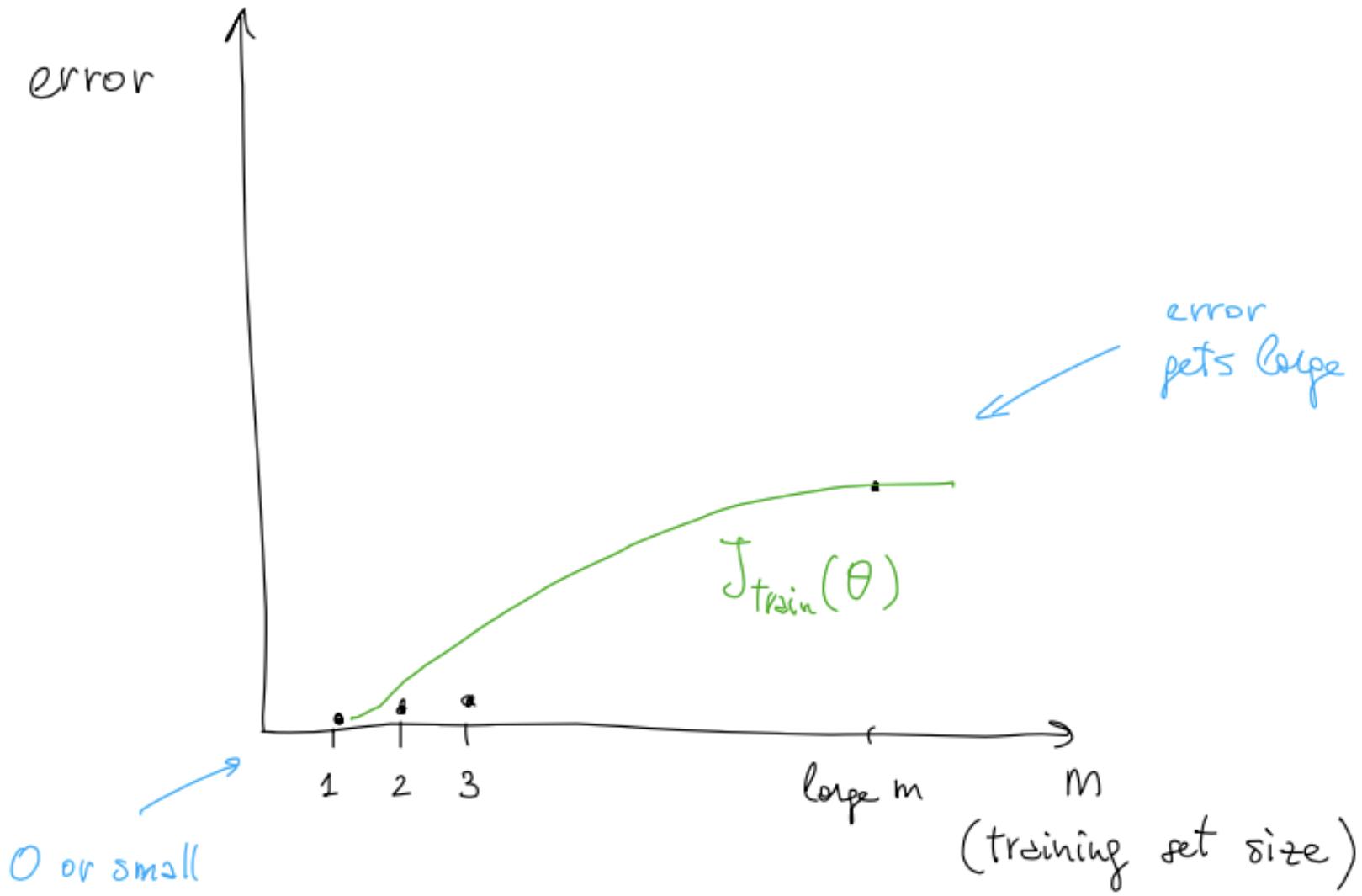
To do so, start with  $m=1$  and add examples one by one ...

Suppose 6 examples, and you try this:

$$h_{\theta}(x) = \theta_0 + \theta_1 x + \theta_2 x^2$$

As  $m \uparrow$ , it becomes harder and harder to fit well all points. So the average training error  $\uparrow$  as  $m \uparrow$ .



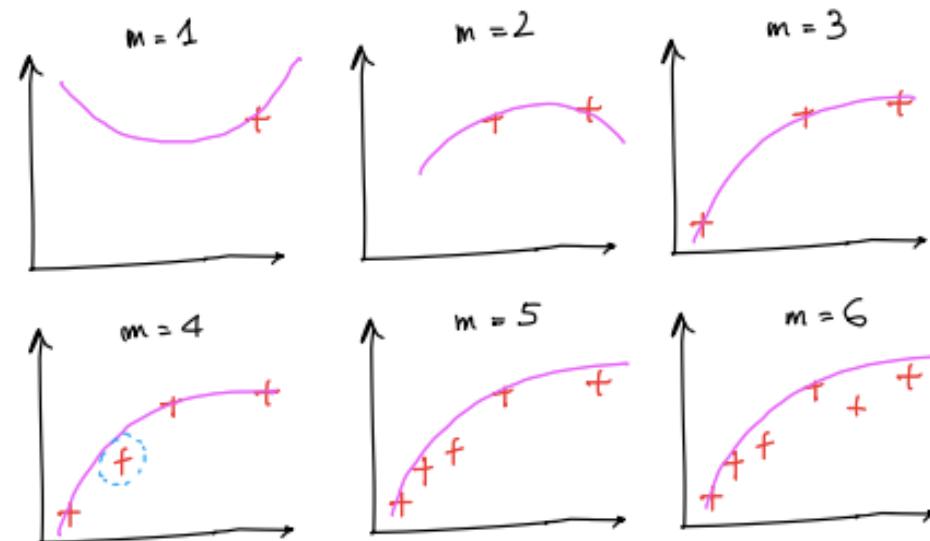


$$J_{\text{train}}(\theta) = \frac{1}{2m_{\text{train}}} \sum_{i=1}^m \left( h_{\theta}(x_{\text{train}}^{(i)}) - y_{\text{train}}^{(i)} \right)^2$$

$$J_{\text{cv}}(\theta) = \frac{1}{2m_{\text{cv}}} \sum_{i=1}^{m_{\text{cv}}} \left( h_{\theta}(x_{\text{cv}}^{(i)}) - y_{\text{cv}}^{(i)} \right)^2$$

Now:

Plot now as a function of  $m$  (training set size)



$$J_{\text{train}}(\theta) = \frac{1}{2m_{\text{train}}} \sum_{i=1}^m (h_{\theta}(x_{\text{train}}^{(i)}) - y_{\text{train}}^{(i)})^2$$

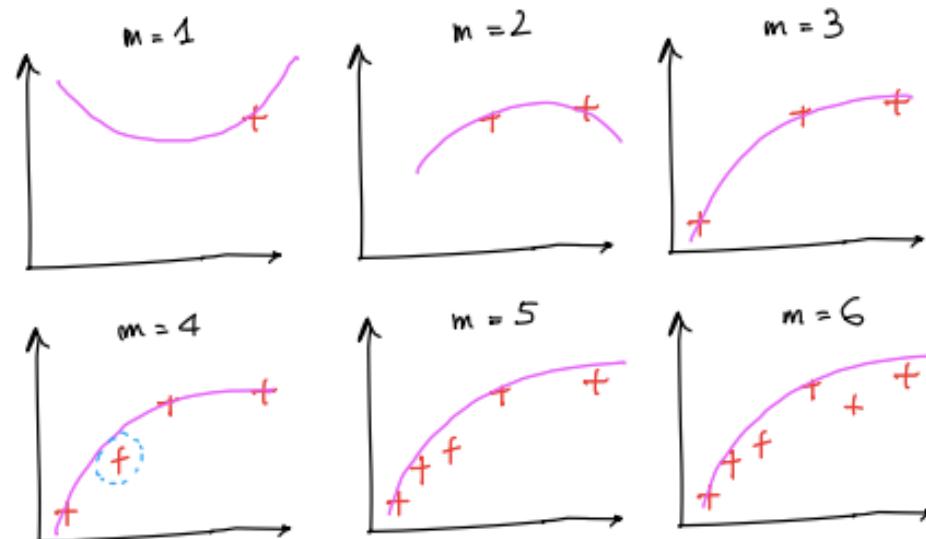
$$J_{\text{cv}}(\theta) = \frac{1}{2m_{\text{cv}}} \sum_{i=1}^{m_{\text{cv}}} (h_{\theta}(x_{\text{cv}}^{(i)}) - y_{\text{cv}}^{(i)})^2$$

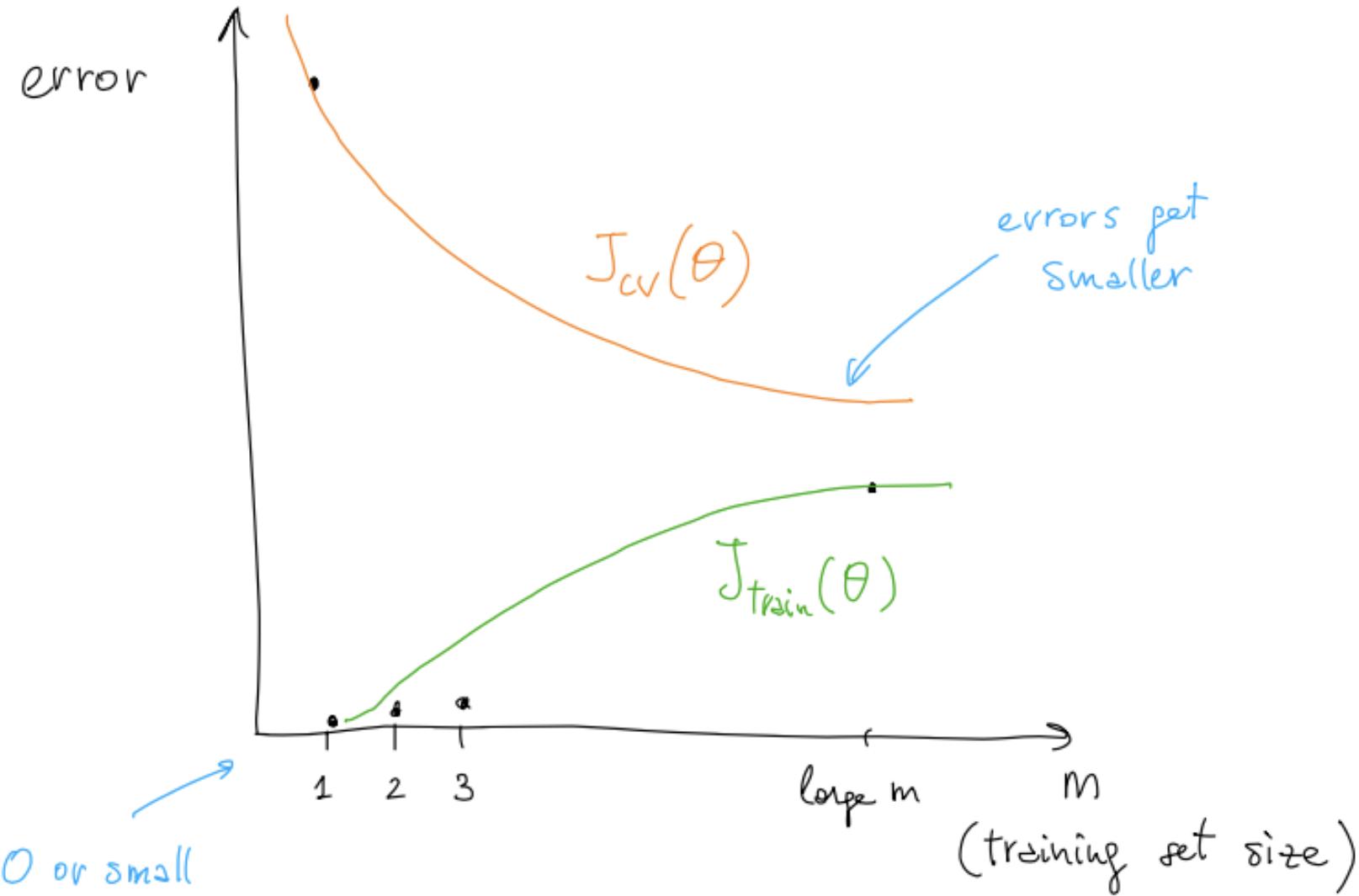
Now:

Plot now as a function of  $m$  (training set size)

$J_{\text{cv}}(\theta)$  tells me how well I can generalize to unseen data  $\Rightarrow$  with small  $m$  values, i.e. a small training set, I cannot generalize well  $\Rightarrow$  CV error would be large. When I get to larger  $m$  I start to see  $h$  fitting data better  $\Rightarrow$  CV error  $\downarrow$

$\downarrow$  the more data you have, the better you generalize to new data.





---

---

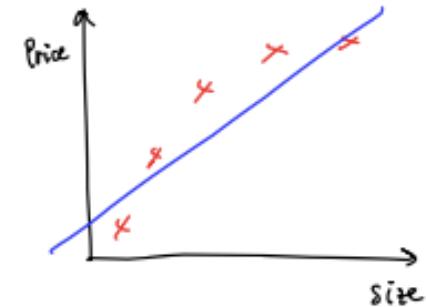
Now let's look at what the learning curves may look like if we have either high bias or high variance problems.

One at a time: **high-bias** first.

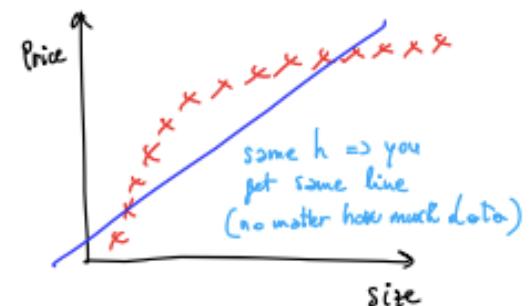
high-bias

→ aim at  $J_{\text{av}}(\theta)$  first

Example:  $h_{\theta}(x) = \Theta_0 + \Theta_1 x$



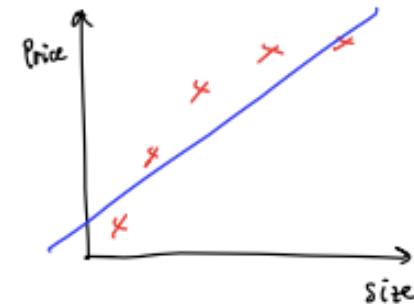
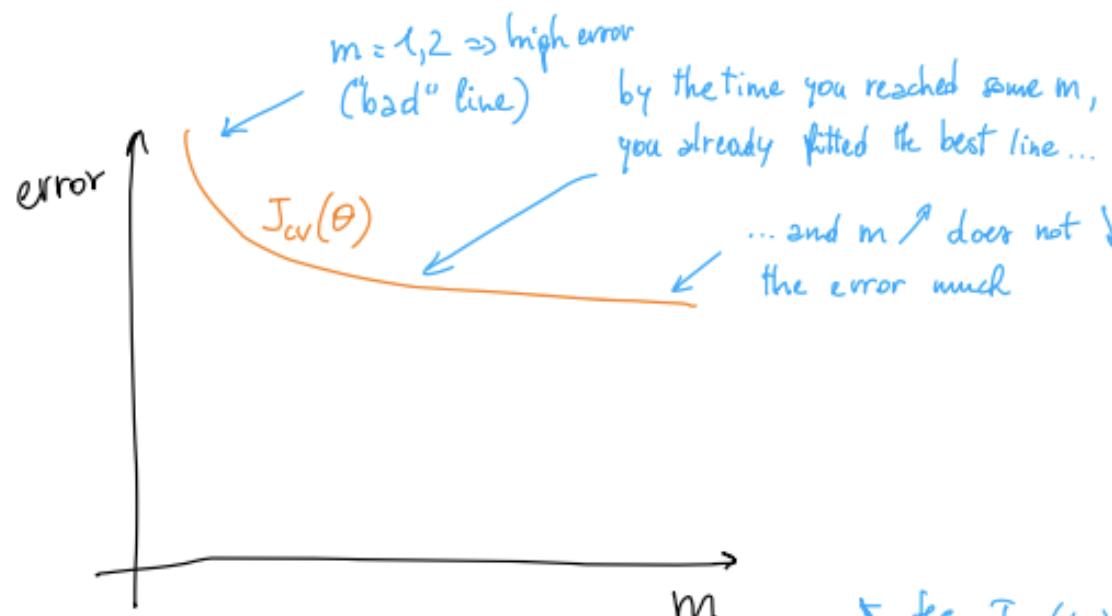
imagine we have  
many more examples



high-bias

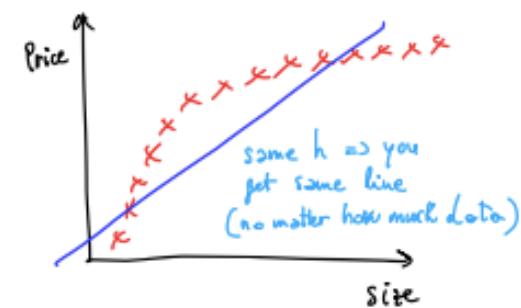
→ aim at  $J_{cv}(\theta)$  first

Example:  $h_{\theta}(x) = \theta_0 + \theta_1 x$

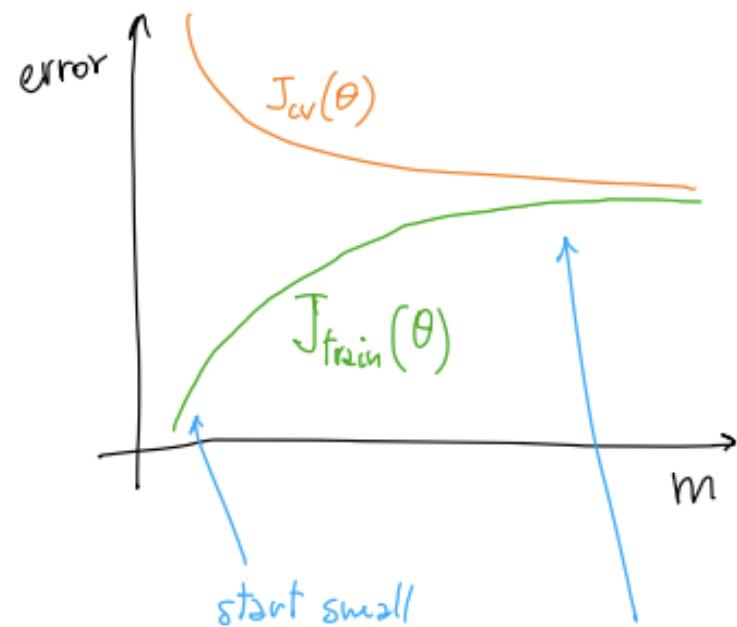


imagine we have  
many more examples

see  $J_{cv}(\theta)$

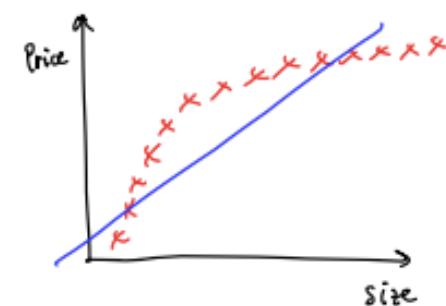
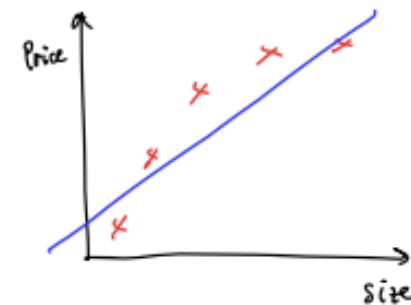


high-bias  $\rightarrow$  sim at  $J_{\text{train}}(\theta)$  now

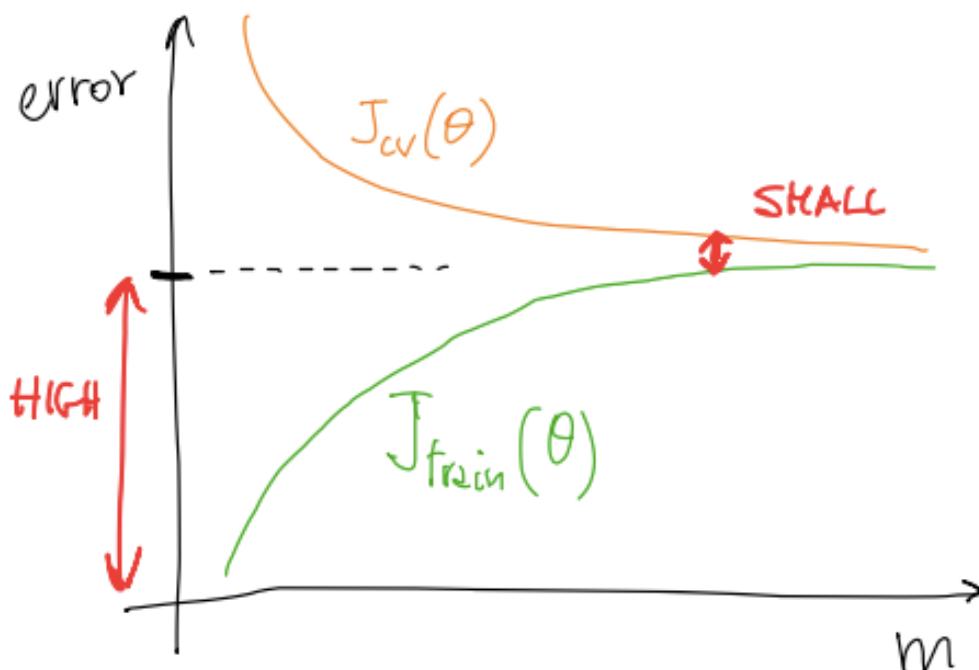


error ↑ and in high-bias  
it ends up close to CV error:  
so few params and so much data  
that the performance on training/cv  
set would be very similar

Example:  $h_{\theta}(x) = \theta_0 + \theta_1 x$



high-bias



The problem with high-bias is that both  $J_{\text{train}}$  and  $J_{\text{cv}}$  are high  $\Rightarrow$  ends up with a relatively high error



IMPORTANT : if a learning algo is suffering from high-bias, getting more training data ( $m \uparrow$ ) will not (by itself) help much



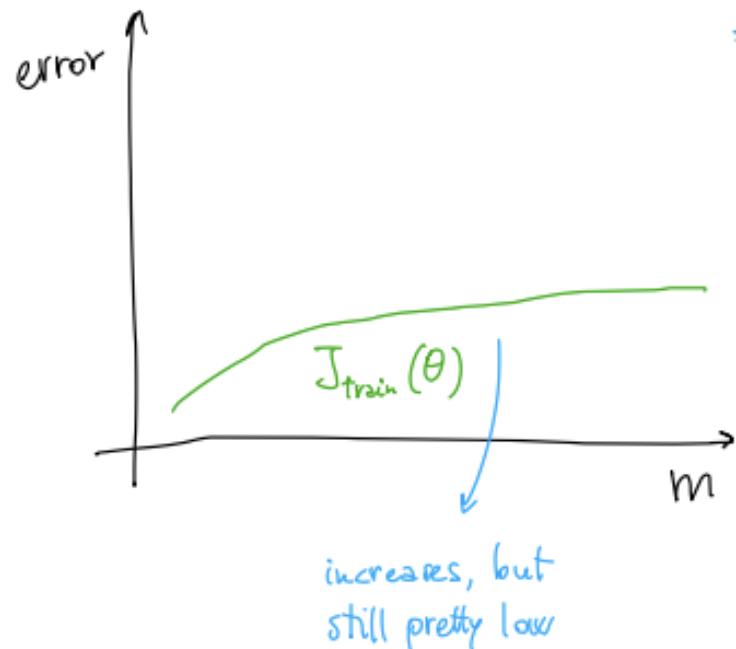
does not help...

---

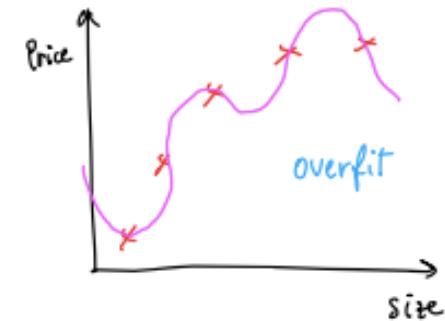
---

And now **high-variance**.

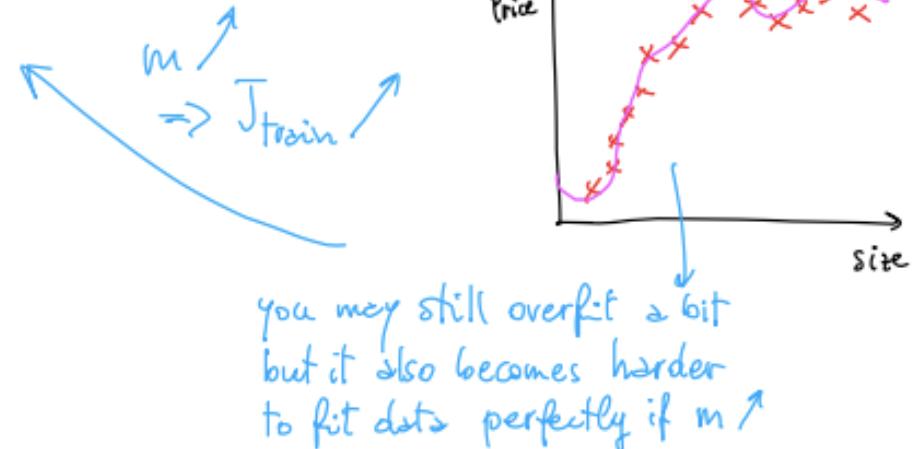
high-variance  $\rightarrow$  small  $J_{\text{train}}(\theta)$  first  $\rightarrow$  Example:  $h_{\theta}(x) = \theta_0 + \theta_1 x + \dots + \theta_{100} x^{100}$   
 (and small  $\lambda$ )



$m$  small  
 $\Rightarrow J_{\text{train}}(\theta)$  small

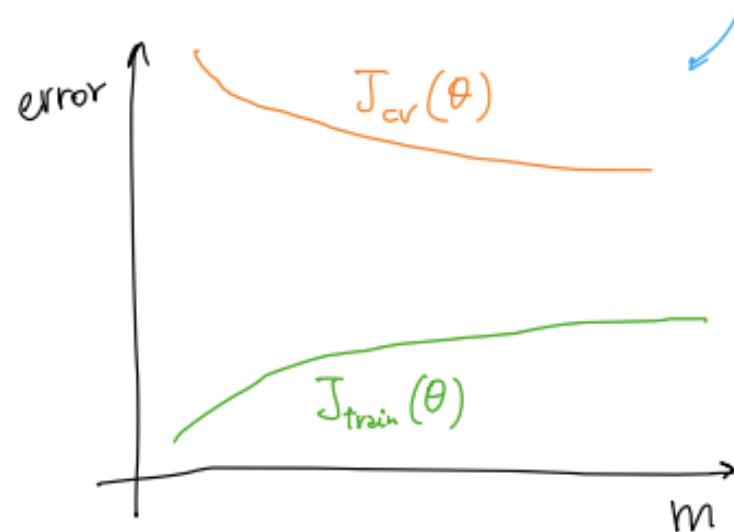


and if  $m$   
increases ...

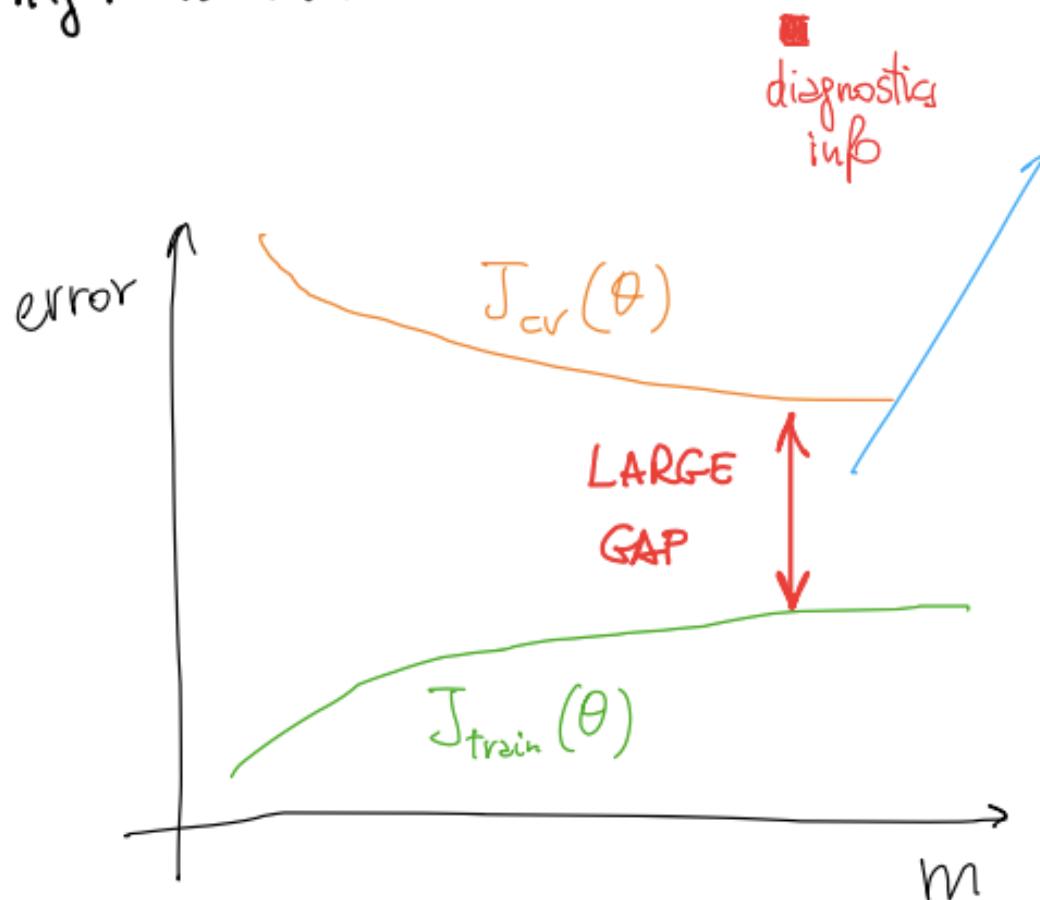


high-variance  $\rightarrow$  sim at  $J_{cv}(\theta)$  now

our  $h$  is overfitting  $\Rightarrow$  cv error would remain high even if we get more examples



high-variance



this large gap is indicative  
(as diagnostic) of a high  
variance scenario.

IMPORTANT : if a learning algo is suffering from high-variance, getting more training data ( $m \uparrow$ ) will likely help !

↓  
extrapolate to  
large  $m$



# Caveats (in all curves so far..)

We have drawn quite clean and idealised curves.

- real life (as said already): more noisy curves, much messier..

Still, plotting learning curves like these can often help you figure out if your learning algorithm is suffering from bias or variance

high-bias



high-variance



# Summary

If I train a learning algo on very few data points (e.g. 2,3) I will easily get error=0 as I can always find e.g. a quadratic function that touches exactly all these data points. So:

- as the training set gets larger, the error for a quadratic function increases
- the error value will plateau out after a certain  $m$  (= training set size)

Algo suffering from **high-bias**:



- small  $m$  (training set size): it causes  $J_{\text{train}}(\Theta)$  to be low and  $J_{\text{cv}}(\Theta)$  to be high
- large  $m$ : it causes both  $J_{\text{train}}(\Theta)$  and  $J_{\text{cv}}(\Theta)$  to be high, with  $J_{\text{train}}(\Theta) \approx J_{\text{cv}}(\Theta)$

⇒ getting more training data will not (by itself!) help much

Algo suffering from **high-variance**:



- small  $m$ : it causes  $J_{\text{train}}(\Theta)$  to be low and  $J_{\text{cv}}(\Theta)$  to be high
- large  $m$ :  $J_{\text{train}}(\Theta)$  increases with  $m$ , and  $J_{\text{cv}}(\Theta)$  decreases with  $m$ , without levelling off. Also,  $J_{\text{train}}(\Theta) < J_{\text{cv}}(\Theta)$  and the difference between them remains significant

⇒ getting more training data is likely to help!

# Applied Machine Learning - Basic

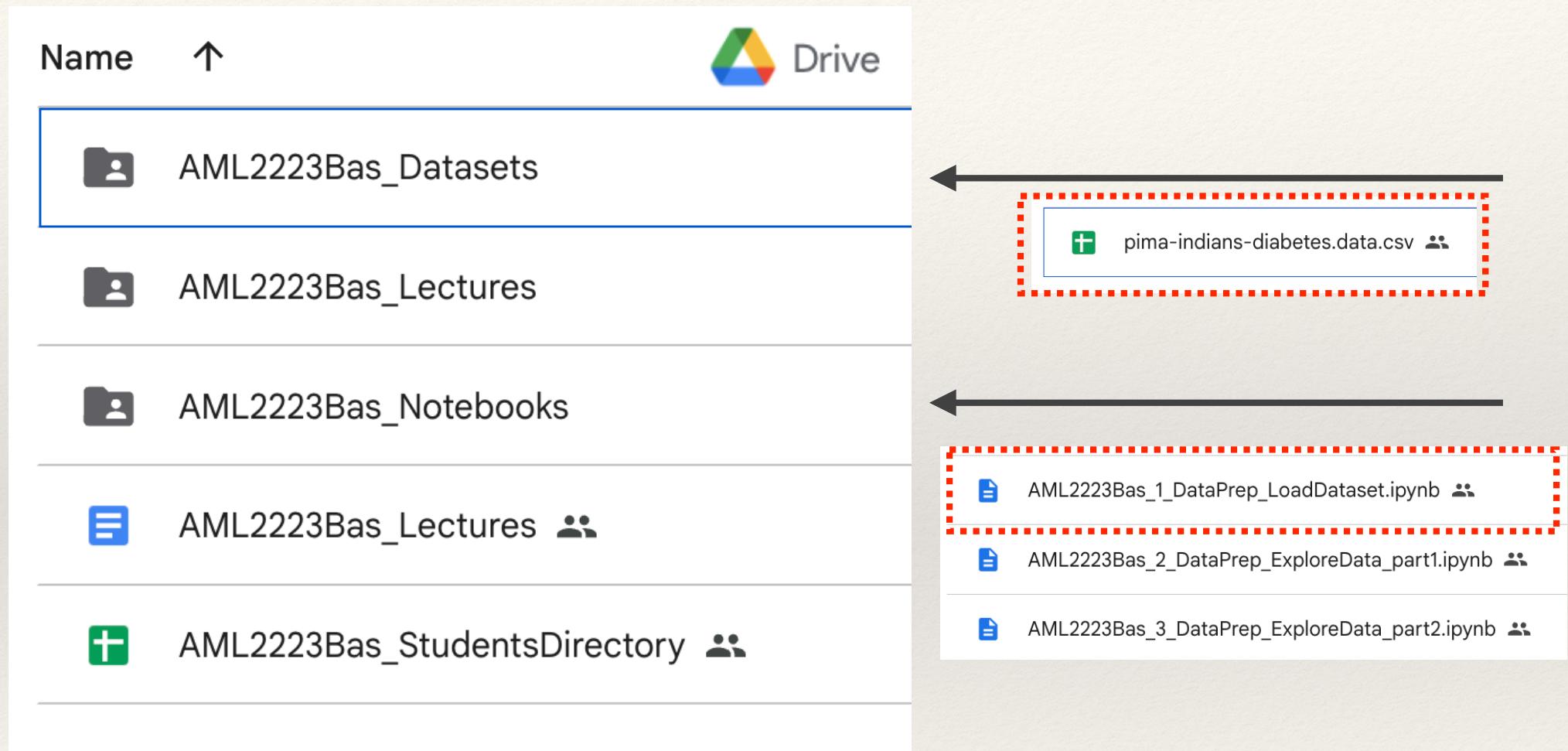
Prof. Daniele Bonacorsi

## Hands-on:

1 - Data prep - load dataset

Data Science and Computation PhD + Master in Bioinformatics  
**University of Bologna**

Look into these:



# Applied Machine Learning - Basic

Prof. Daniele Bonacorsi

## Hands-on:

### 2 - Data prep - explore data (part 1)

Data Science and Computation PhD + Master in Bioinformatics  
**University of Bologna**

Look into these:

The screenshot shows a Google Drive interface with a list of items:

- Name ↑
- Drive
- AML2223Bas\_Datasets
- AML2223Bas\_Lectures
- AML2223Bas\_Notebooks
- AML2223Bas\_Lectures 2
- AML2223Bas\_StudentsDirectory 2

A red dashed box highlights the following items:

- pima-indians-diabetes.data.csv
- AML2223Bas\_1\_DataPrep\_LoadDataset.ipynb
- AML2223Bas\_2\_DataPrep\_ExploreData\_part1.ipynb
- AML2223Bas\_3\_DataPrep\_ExploreData\_part2.ipynb