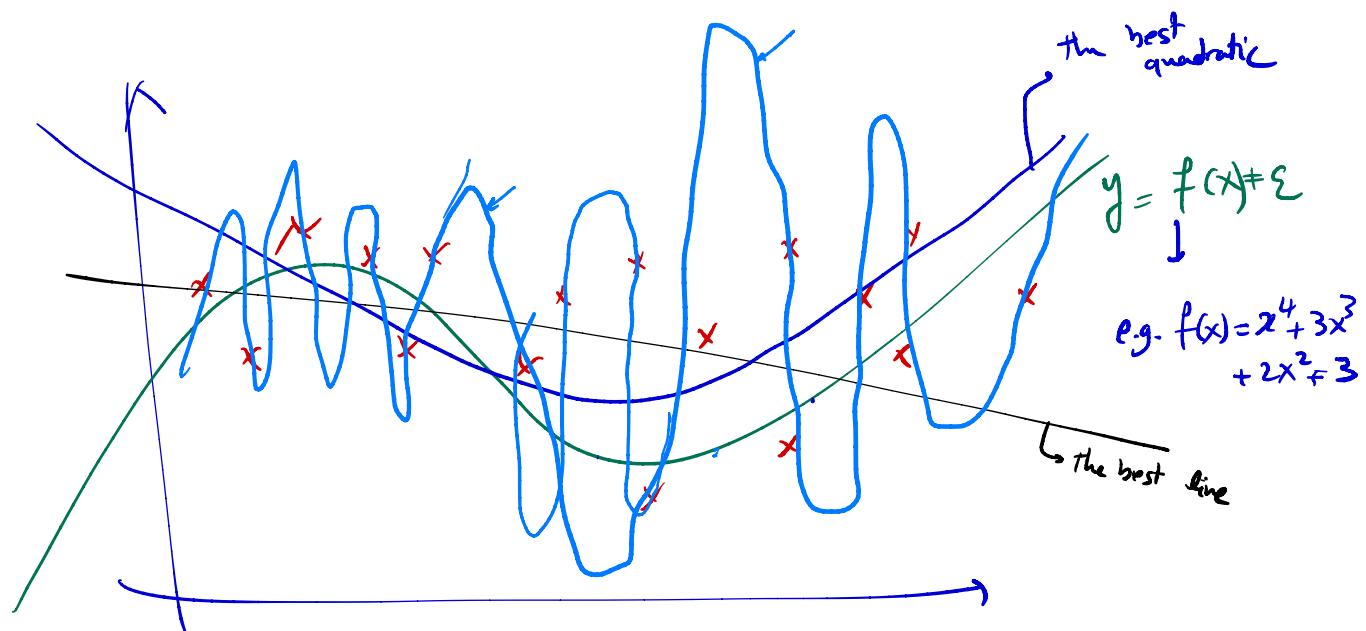
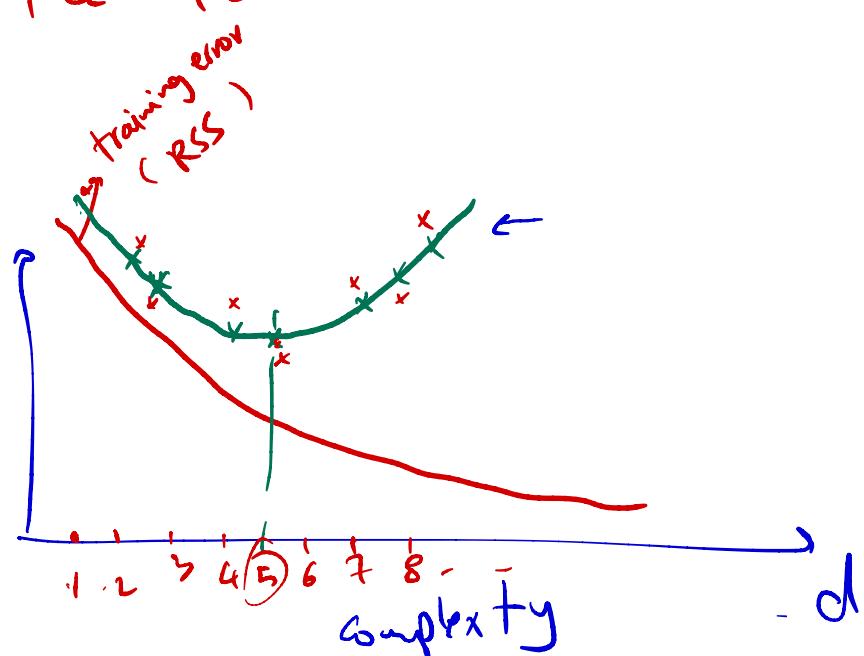


## Lecture 16 :



the best  
fit / model with  
degree 100

Question : How can we estimate the  
the test error / curve ?

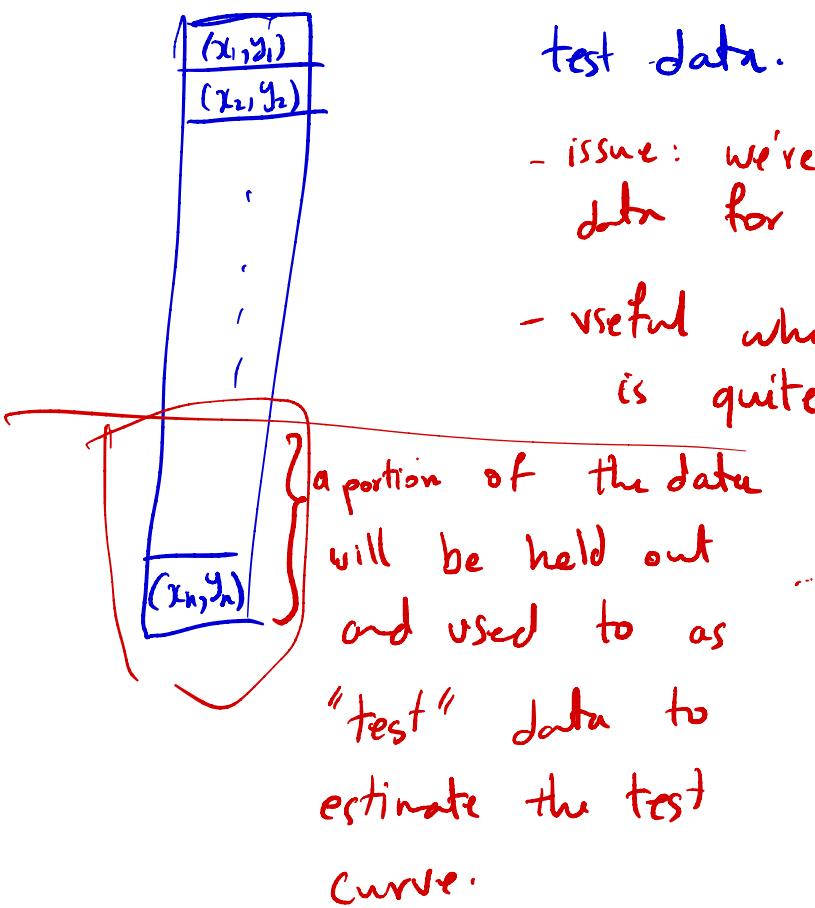


Two methods / answers :

(1) validation / hold-out data:

Randomly choose a (given) portion of the data and use it as test data. Use the rest for training.

- issue: we're not using all the available data for training
- useful when the size of the data is quite large



(2) k-fold cross-validation:

partition / Divide the data randomly into  $k$  parts:

$$D = D_1 \cup D_2 \cup D_3 \dots \cup D_k$$

for each  $i=1, \dots, k$

- train on  $D \setminus D_i$
- test the resulting model on  $D_i$

$\Rightarrow \text{test}_i = \text{test error of the trained model on } D_i$

$$\text{test-error} \approx \frac{1}{k} (\text{test}_1 + \text{test}_2 + \dots + \text{test}_k)$$

---

$\Rightarrow$  both of the methods are used to estimate the test curve

- once the "best complexity" is found (minimum in the estimated test curve),  $d^{\text{best}}$  we will use the whole data set to train a model with the "best complexity"
- cross validation is done for each complexity to estimate the test error.



# Classification:

Regression :

label  $y_i \in \mathbb{R}$

Classification:

label  $\in$  Discrete set

$y_i \in \{0, 1\}$

$\in \{\text{Spam, not-spam}\}$

$\in \{\text{Sunny, Snow, rain}\}$

Example:

$$\{(x_i, y_i)\}_{i=1}^n$$

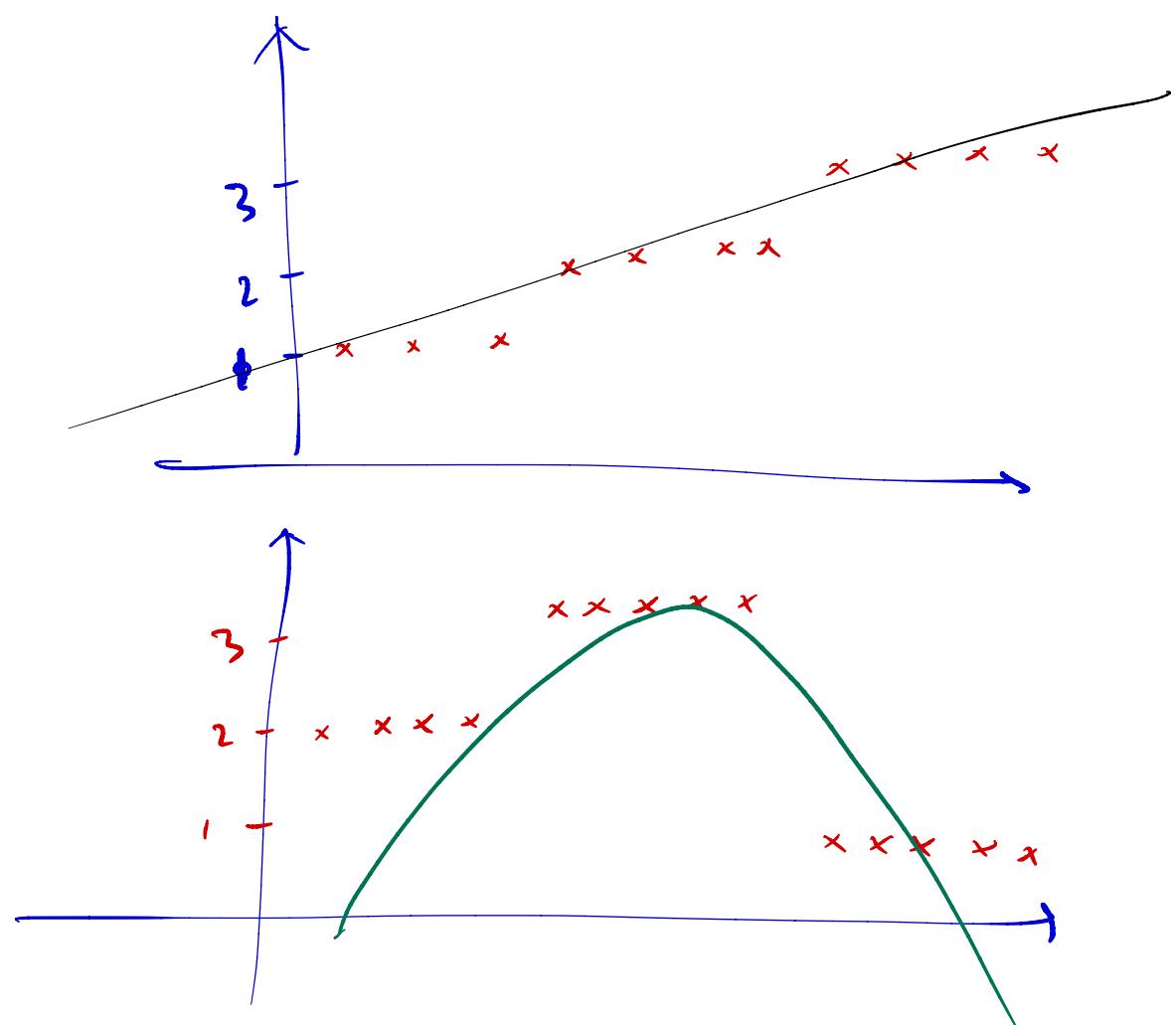
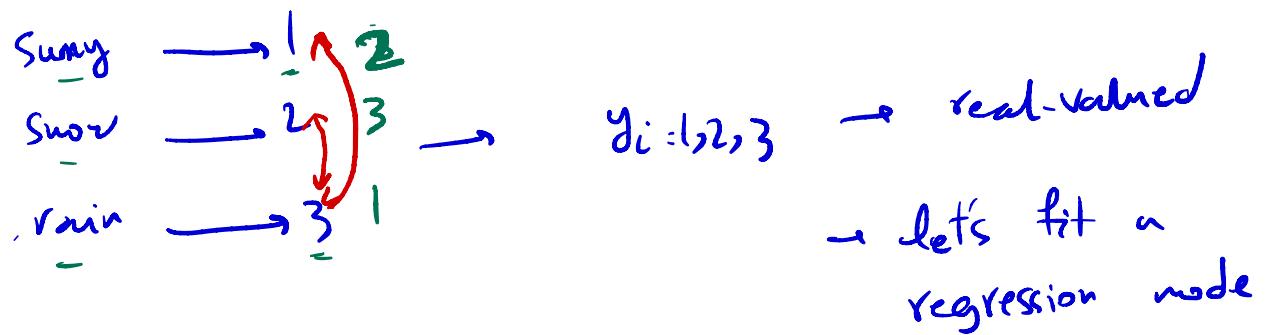
$x_i$  = weather data for day  $i$  ((clouds, --))

$y_i = \{\text{sunny, snowy, rainy}\}$

$$\underline{y_i} \approx f(\underline{x_i})$$

(remember, we want to be accurate at test time.)

As a naive approach, let's try model/solve  
a classification problem using regression:



Example:

input  $x$  (input data)  
weather data

$$\left. \begin{array}{l} \Pr\{\text{Sunny} \mid x\} = 0.2 \\ \Pr\{\text{Snow} \mid x\} = 0.5 \\ \Pr\{\text{Rain} \mid x\} = 0.3 \end{array} \right\}$$

- what we really need to learn/estimate  
is the "conditional probabilities":

$$\Pr\{\text{class}_k \mid x\} \quad \text{for } k=1,\dots,K$$

Bayes Optimal Classifier:

---

Assume for the moment that all the  
conditional probabilities are known

$$\Pr\{\text{class}_k \mid x\} \quad \text{known for } k=1,\dots,K.$$

How should we (optimally) predict the

class associated to input  $x$ ?

Bayes optimal classifier:

$$\begin{aligned}\hat{y}_{\text{Bayes}}(x) &= \text{predicted class}(x) \\ &= \underset{k}{\operatorname{argmax}} \Pr\{\text{class}_k | x\}\end{aligned}$$

(break ties arbitrarily)

Theorem: The Bayes optimal classifier has the smallest classification error among all the possible classifiers.

This means: for every other classifier

$\hat{y}_{\text{other}}(x)$ , we have:

$$\Pr_{(x,y)}\{\hat{y}_{\text{other}}(x) \neq y\} \geq \Pr_{(x,y)}\{\hat{y}_{\text{Bayes}}(x) \neq y\}$$

- The fundamental problem with the Bayes classifier is that it assumes the knowledge of the conditional probabilities (which not true in practice)
- All the methods for classification in ML aim at estimating the conditional Probabilities.
- We'll cover several methods to estimate these conditional probabilities from data.

## Logistic Regression:

For now, assume that we have two classes (the binary classification problem):

Class<sub>0</sub>  $\rightarrow y_i = 0$

$x_i \in \mathbb{R}$

Class<sub>1</sub>  $\rightarrow y_i = 1$

Goal: given  $x \rightarrow \underbrace{\left\{ \begin{array}{l} \Pr\{0|x\} \\ \Pr\{1|x\} \end{array} \right\}}$

Data:  $\left\{ (x_i, y_i) \right\}_{i=1}^n$

Step 1:  
parametric model

↓  
Step 2:  
fit the  
params to data

Step: we're going to design a parametric model for estimation

$$\underbrace{\Pr\{0|x\}}_{\in [0,1]}, \quad \underbrace{\Pr\{1|x\}}_{\in [0,1]}$$

Let's begin with linear models.

What is a linear model in the binary classification setting?

$$\Pr\{0|x\} = f(x; \underline{\text{parameters}})$$

$$\Pr\{1|x\} = 1 - f(x; \underline{\text{parameters}})$$

let's come up with simplest candidates for the parametric model.

- Perhaps the simplest model would be something like

$$\begin{cases} \Pr\{0|x\} = \beta_0 + \beta_1 x \\ \Pr\{1|x\} = 1 - (\beta_0 + \beta_1 x) \end{cases}$$

But, one can not enforce the constraint that the probabilities are between 0 and 1 using the above model (i.e. it's not possible to ensure that  $\beta_0 + \beta_1 x \in [0, 1]$ ).

Instead, we consider the following model:

$$\Pr\{0|x\} = \frac{e^{B_0 + B_1 x}}{1 + e^{B_0 + B_1 x}}$$

$$\Pr\{1|x\} = 1 - \Pr\{0|x\}$$

$$= \frac{1}{1 + e^{B_0 + B_1 x}}$$

It's easy to see that  $\frac{e^{B_0 + B_1 x}}{1 + e^{B_0 + B_1 x}} \in [0,1]$