

# EC-350 AI and Decision Support Systems

## Week 9 K-Nearest Neighbour Classifier

Dr. Arslan Shaukat



Acknowledgements: Lecture slides material from  
Duda, Hart and Stork, Dr. Gavin Brown

### Problem Statement

- Why recognising rugby players is (almost) the same problem as *recognising handwritten digits*



7210414959  
0690159784  
9665407401  
3134727121  
1742351244



## Problem Statement

Can we LEARN to recognise a rugby player and ballet dancer?



What are the “features” of a rugby player?

06/12/2017

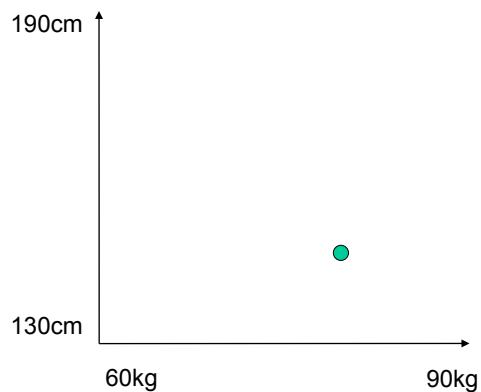
EC-350 AI and DSS

EME (NUST)

3

## Features

Rugby players = short + heavy?



06/12/2017

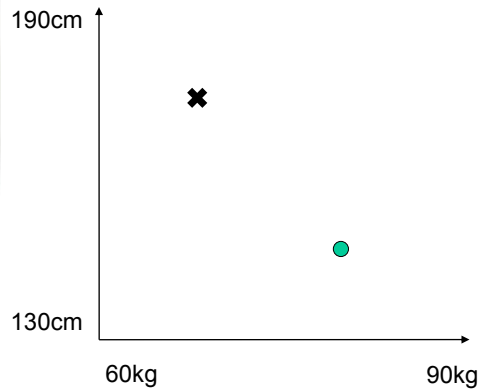
EC-350 AI and DSS

EME (NUST)

4

## Features

Ballet dancers = tall + skinny?



06/12/2017

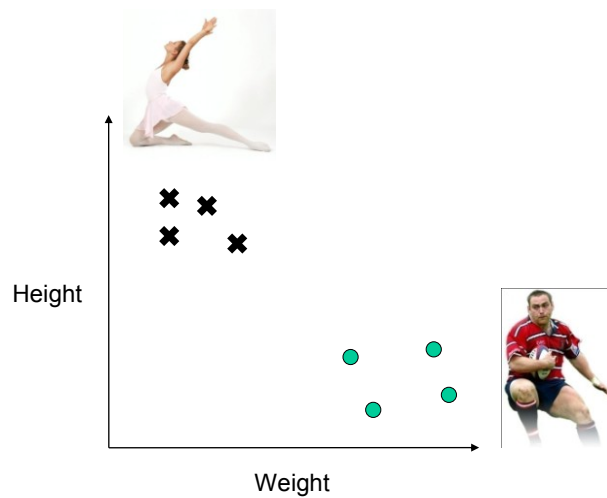
EC-350 AI and DSS

EME (NUST)

5

## Feature Space

Rugby players “cluster” separately in the space.



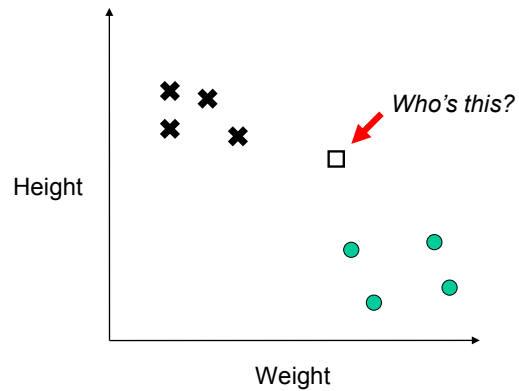
06/12/2017

EC-350 AI and DSS

EME (NUST)

6

## The K-Nearest Neighbour Algorithm



06/12/2017

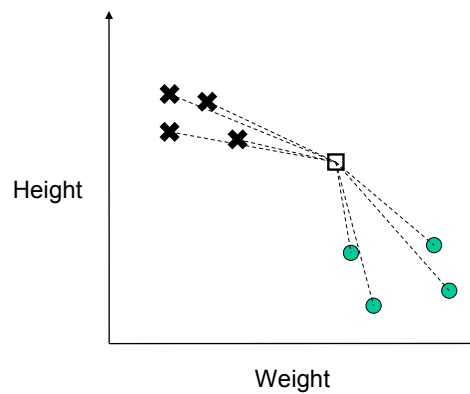
EC-350 AI and DSS

EME (NUST)

7

## The K-Nearest Neighbour Algorithm

1. Measure distance to all points



06/12/2017

EC-350 AI and DSS

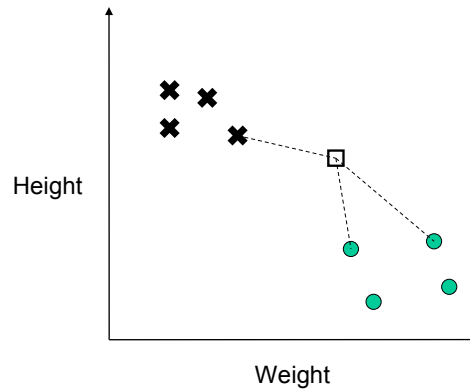
EME (NUST)

8

## The K-Nearest Neighbour Algorithm

1. Measure distance to all points
2. Find closest "k" points

← (here k=3, but it could be more)



06/12/2017

EC-350 AI and DSS

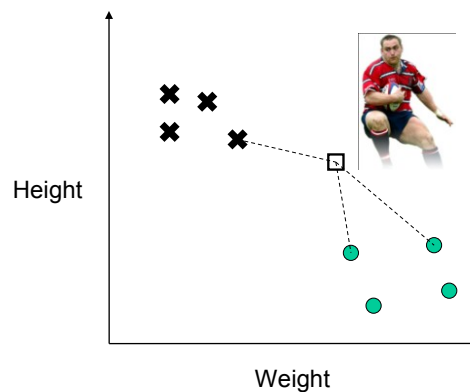
EME (NUST)

9

## The K-Nearest Neighbour Algorithm

1. Measure distance to all points
2. Find closest "k" points
3. Assign majority class

← (here k=3, but it could be more)



06/12/2017

EC-350 AI and DSS

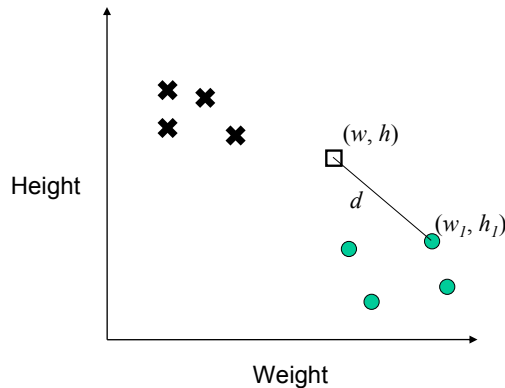
EME (NUST)

10

## Distance Measure

“Euclidean distance”

$$d = \sqrt{(w - w_1)^2 + (h - h_1)^2}$$



06/12/2017

EC-350 AI and DSS

EME (NUST)

11

## The K-Nearest Neighbour Algorithm

for each testing point

*measure distance to every training point*

*find the  $k$  closest points*

*identify the most common class among those  $k$*

*assign that class*

end

- Advantage: Surprisingly good classifier!
- Disadvantage: Have to store the entire training set in memory

06/12/2017

EC-350 AI and DSS

EME (NUST)

12

## Distance Measure

Euclidean distance still works in 3-d, 4-d, 5-d, etc....

$$d = \sqrt{(x - x_1)^2 + (y - y_1)^2 + (z - z_1)^2}$$

$x = \text{Height}$ $y = \text{Weight}$ $z = \text{Shoe size}$
--

06/12/2017

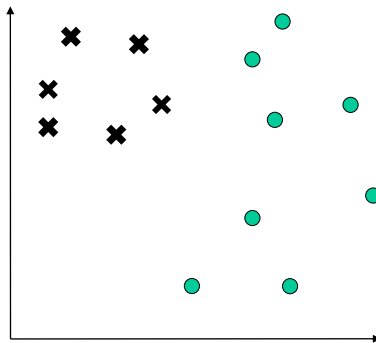
EC-350 AI and DSS

EME (NUST)

13

## Over-fitting

An **Important** Concept in Machine Learning



06/12/2017

EC-350 AI and DSS

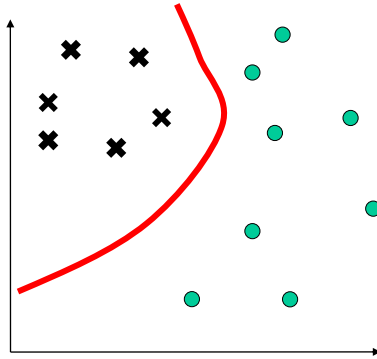
EME (NUST)

14

## Over-fitting

An **Important** Concept in Machine Learning

*Looks good so far...*



06/12/2017

EC-350 AI and DSS

EME (NUST)

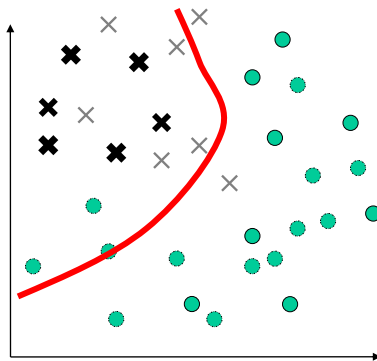
15

## Over-fitting

An **Important** Concept in Machine Learning

*Looks good so far...*

*Oh no! Mistakes!  
What happened?*



06/12/2017

EC-350 AI and DSS

EME (NUST)

16

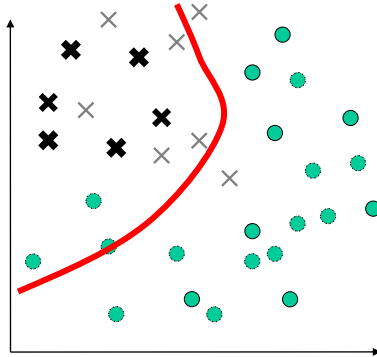


## Over-fitting

### An Important Concept in Machine Learning

*Looks good so far...*

*Oh no! Mistakes!  
What happened?*



We didn't have all the data.

We can never assume that we do.

This is called "OVER-FITTING"  
to the small dataset.

06/12/2017

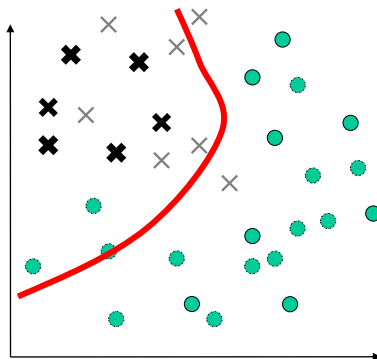
EC-350 AI and DSS

EME (NUST)

17

## Over-fitting

- While an overly complex boundary may allow perfect classification of the training samples, it is unlikely to give good classification of novel patterns



06/12/2017

EC-350 AI and DSS

EME (NUST)

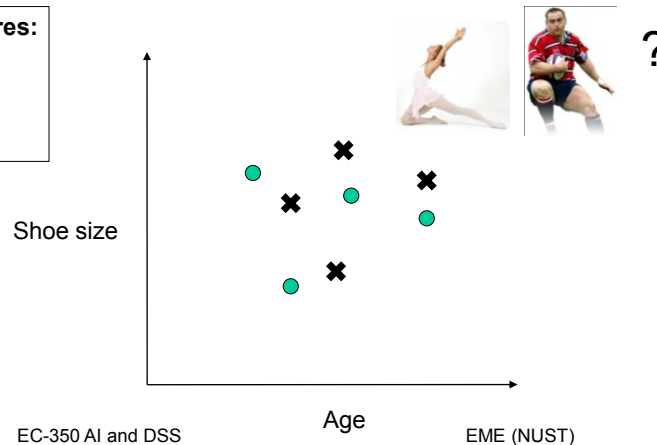
18

## Features

Choosing the wrong features makes it difficult  
Too many features  
It's computationally intensive

### Possible features:

- Shoe size ✓
- Height
- Age ✓
- Weight



06/12/2017

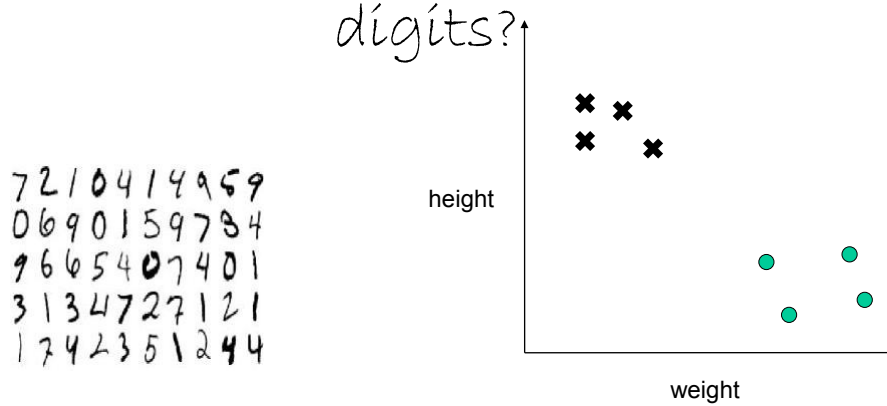
EC-350 AI and DSS

EME (NUST)

19

## Our Problem

Now – how is this problem like  
recognising handwritten  
digits?



06/12/2017

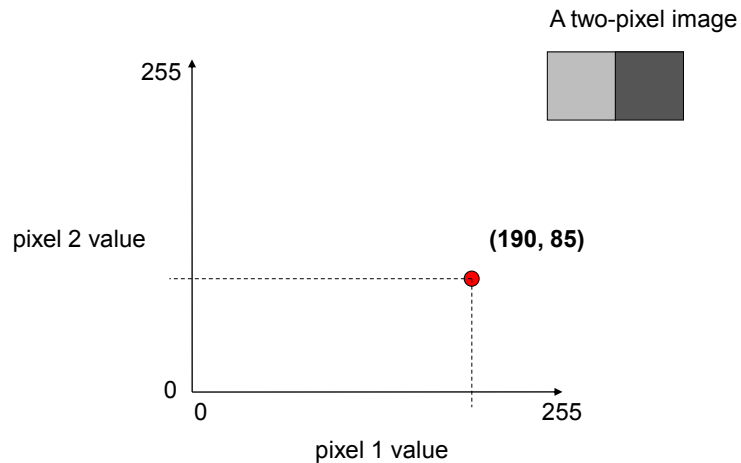
EC-350 AI and DSS

EME (NUST)

20

## Recognising Handwritten Digits

Let's say the axes now represent **pixel values**.



06/12/2017

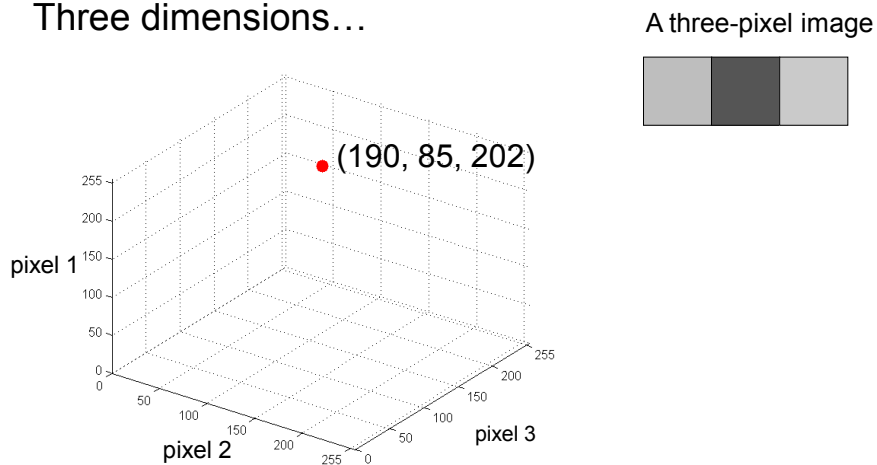
EC-350 AI and DSS

EME (NUST)

21

## Recognising Handwritten Digits

Three dimensions...



This 3-pixel image is represented by a **SINGLE** point in a 3-D space.

06/12/2017

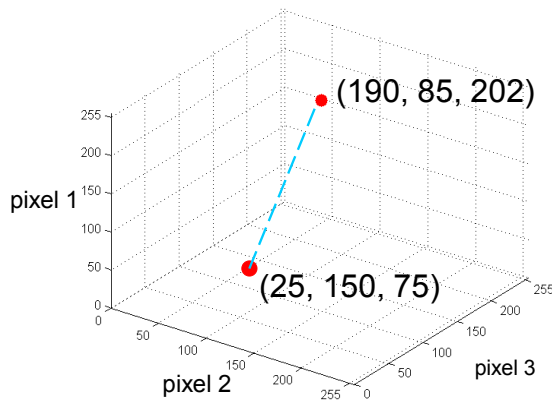
EC-350 AI and DSS

EME (NUST)

22

## Recognising Handwritten Digits

### Distances between images



A three-pixel image



Another 3-pixel image



Straight line distance between them

06/12/2017

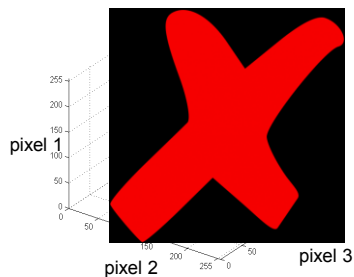
EC-350 AI and DSS

EME (NUST)

23

## Recognising Handwritten Digits

### Four dimensions? Five? Six-dimensional spaces?



A three-pixel image



A four-pixel image.



A five-pixel image



06/12/2017

EC-350 AI and DSS

EME (NUST)

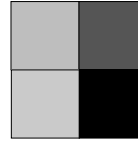
24

## Recognising Handwritten Digits

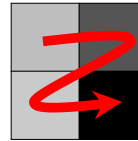
A four-pixel image.



A different four-pixel image.



(190, 85, 202, 10)  
Same 4-dimensional vector!



Assuming we read pixels in a systematic manner, we can now represent any image as a single point in a high dimensional space.

06/12/2017

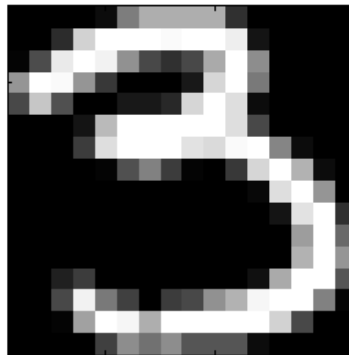
EC-350 AI and DSS

EME (NUST)

25

## Recognising Handwritten Digits

16 x 16 image.... How many dimensions?



06/12/2017

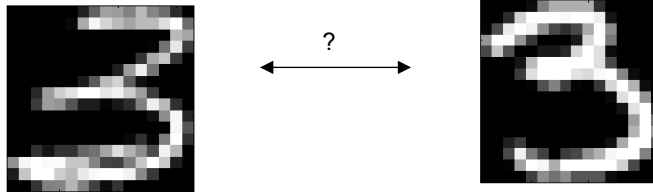
EC-350 AI and DSS

EME (NUST)

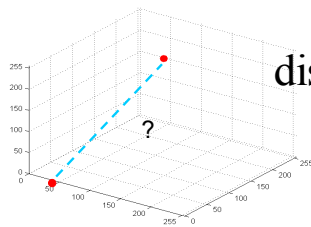
26

## Recognising Handwritten Digits

Distances between digits....



Straight-line distance in N-dimensional space?



$$\text{distance}(x, x') = \sqrt{\sum_{i=1}^{i=256} (x_i - x'_i)^2}$$

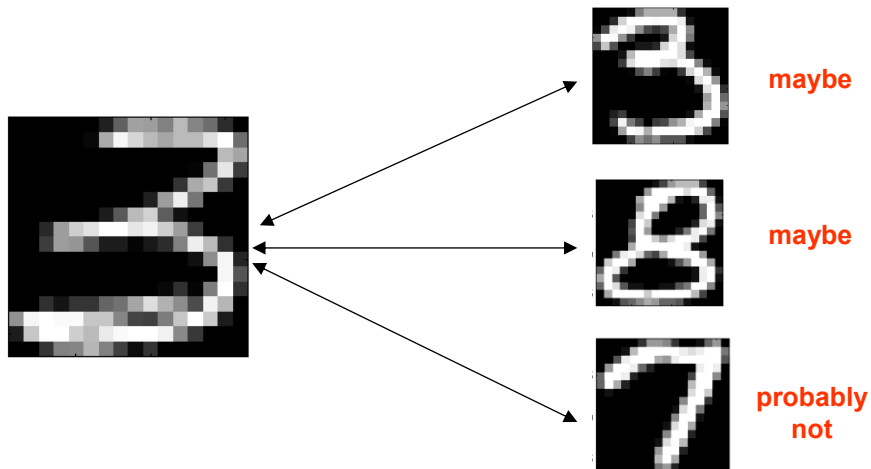
06/12/2017

EC-350 AI and DSS

EME (NUST)

27

## Recognising Handwritten Digits



Which is **closest neighbour** in N-dimensions?

06/12/2017

EC-350 AI and DSS

EME (NUST)

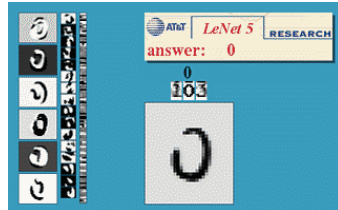
28

# Recognising Handwritten Digits

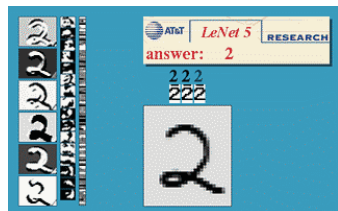


## AT&T Research Labs

The USPS ZIP code reader – **learnt** from examples.



FAST recognition



NOISE resistant, can **generalise** to future unseen patterns

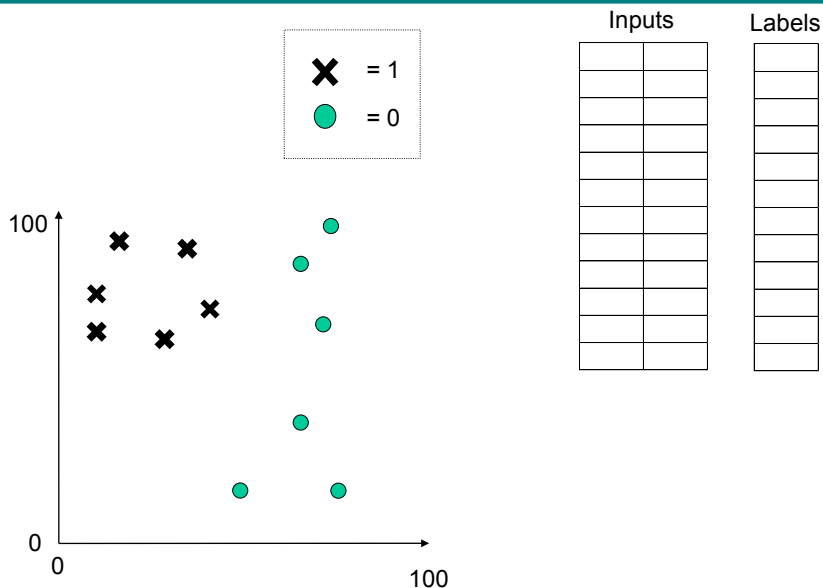
06/12/2017

EC-350 AI and DSS

EME (NUST)

29

## Training and Testing Data



06/12/2017

EC-350 AI and DSS

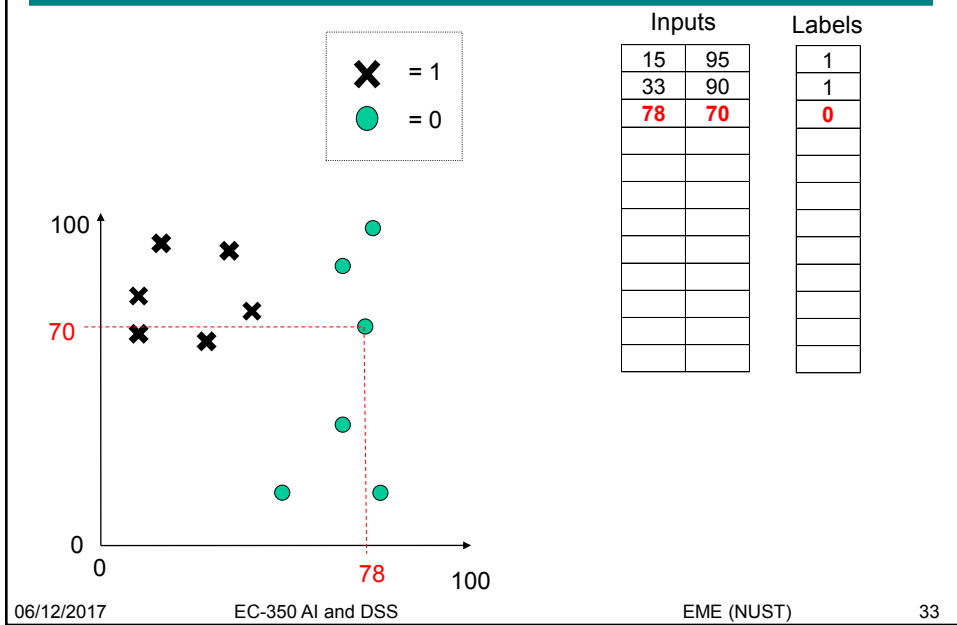
EME (NUST)

30

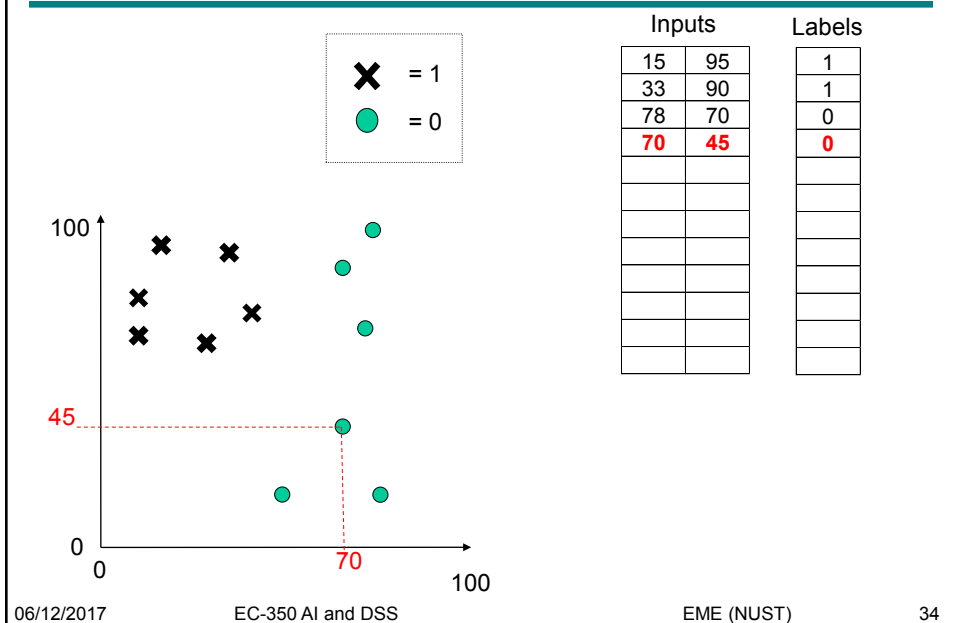




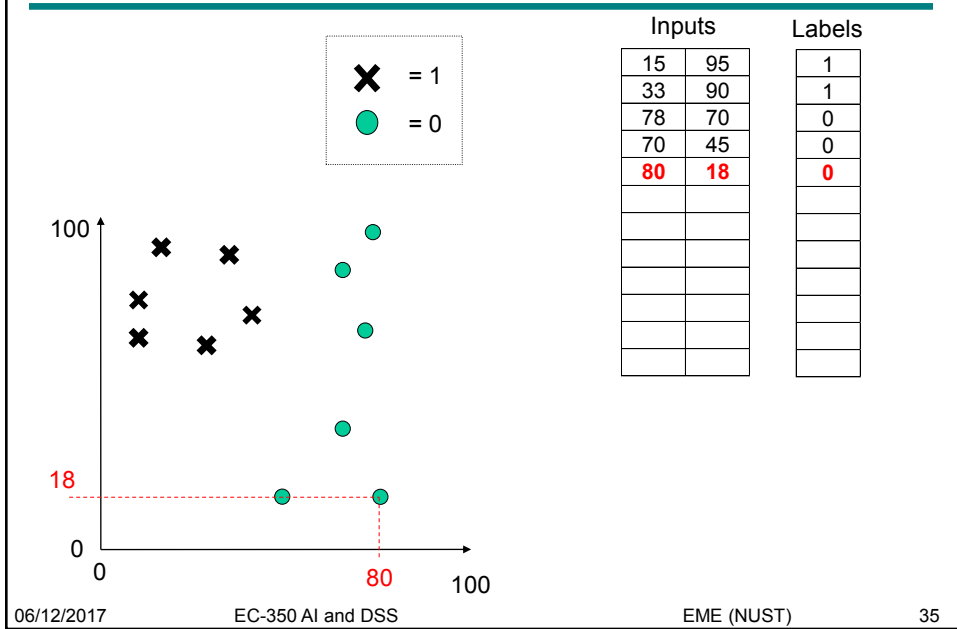
## Training and Testing Data



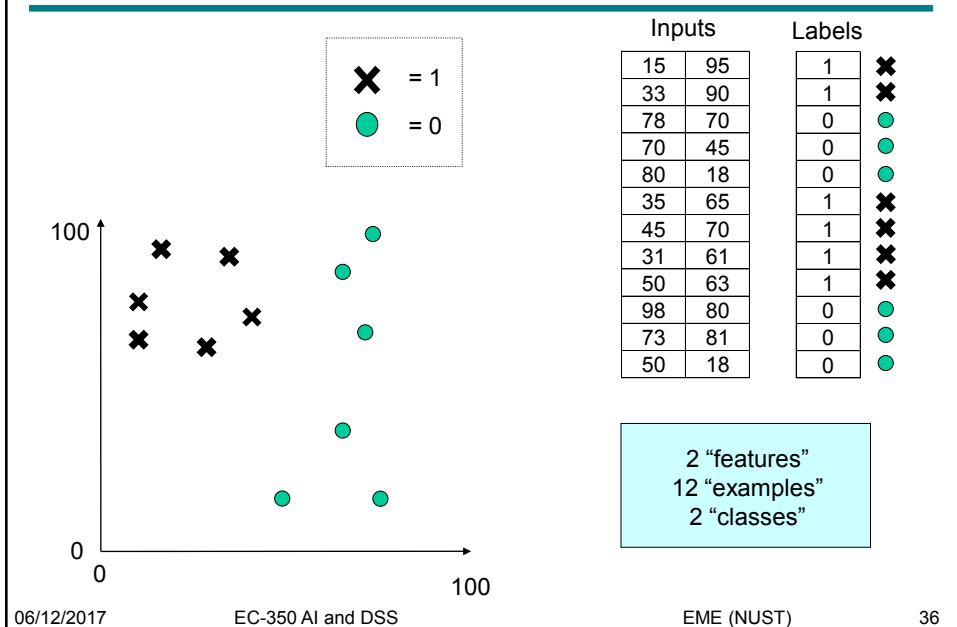
## Training and Testing Data



## Training and Testing Data



## Training and Testing Data



## Training and Testing Data

### Dataset

Inputs		Labels
15	95	1
33	90	1
78	70	0
70	45	0
80	18	0
35	65	1
45	70	1
31	61	1
50	63	1
98	80	0
73	81	0
50	18	0



06/12/2017

EC-350 AI and DSS

EME (NUST)

37

## Training and Testing Data

Inputs		Labels																																							
15	95	1	 50:50 split 	<table><tr><td>15</td><td>95</td><td>1</td></tr><tr><td>33</td><td>90</td><td>1</td></tr><tr><td>78</td><td>70</td><td>0</td></tr><tr><td>70</td><td>45</td><td>0</td></tr><tr><td>80</td><td>18</td><td>0</td></tr><tr><td>35</td><td>65</td><td>1</td></tr></table>	15	95	1	33	90	1	78	70	0	70	45	0	80	18	0	35	65	1	<table><tr><td>15</td><td>95</td><td>1</td></tr><tr><td>33</td><td>90</td><td>1</td></tr><tr><td>78</td><td>70</td><td>0</td></tr><tr><td>70</td><td>45</td><td>0</td></tr><tr><td>80</td><td>18</td><td>0</td></tr><tr><td>35</td><td>65</td><td>1</td></tr></table>	15	95	1	33	90	1	78	70	0	70	45	0	80	18	0	35	65	1
15	95	1																																							
33	90	1																																							
78	70	0																																							
70	45	0																																							
80	18	0																																							
35	65	1																																							
15	95	1																																							
33	90	1																																							
78	70	0																																							
70	45	0																																							
80	18	0																																							
35	65	1																																							
45	70	1	<table><tr><td>45</td><td>70</td><td>1</td></tr><tr><td>31</td><td>61</td><td>1</td></tr><tr><td>50</td><td>63</td><td>1</td></tr><tr><td>98</td><td>80</td><td>0</td></tr><tr><td>73</td><td>81</td><td>0</td></tr><tr><td>50</td><td>18</td><td>0</td></tr></table>	45	70	1	31	61	1	50	63	1	98	80	0	73	81	0	50	18	0	<table><tr><td>45</td><td>70</td><td>1</td></tr><tr><td>31</td><td>61</td><td>1</td></tr><tr><td>50</td><td>63</td><td>1</td></tr><tr><td>98</td><td>80</td><td>0</td></tr><tr><td>73</td><td>81</td><td>0</td></tr><tr><td>50</td><td>18</td><td>0</td></tr></table>	45	70	1	31	61	1	50	63	1	98	80	0	73	81	0	50	18	0	
45	70	1																																							
31	61	1																																							
50	63	1																																							
98	80	0																																							
73	81	0																																							
50	18	0																																							
45	70	1																																							
31	61	1																																							
50	63	1																																							
98	80	0																																							
73	81	0																																							
50	18	0																																							
31	61	1																																							
50	63	1																																							
98	80	0																																							
73	81	0																																							
50	18	0																																							

06/12/2017

EC-350 AI and DSS

EME (NUST)

38

## Training and Testing Data

*If too small... you'll make many mistakes on the testing set.*

*Needs to be big.  
Bigger is better.*

15	95	1
33	90	1
78	70	0
70	45	0
80	18	0
35	65	1

### Training set

*Train a K-NN on this...*

45	70	1
31	61	1
50	63	1
98	80	0
73	81	0
50	18	0

### Testing set

*... then, test it on this!*

*"simulates" what it might be like to see new data in the future*

06/12/2017

EC-350 AI and DSS

EME (NUST)

39

## Training and Testing Data

### Training set

Build a k-NN using training set

### Testing set

Percentage of incorrect predictions is called the "error"

e.g. "Training" error  
e.g. "Testing" error

How many incorrect predictions on testing set?

06/12/2017

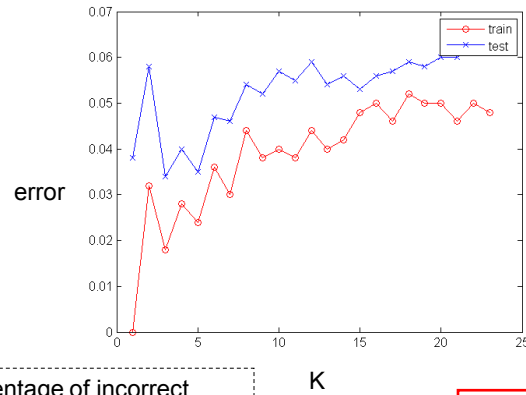
EC-350 AI and DSS

EME (NUST)

40

## Classifying '3' versus '8' Digits

Plotting error as a function of 'K' (as in the K-NN)



Percentage of incorrect predictions is called the "error"

e.g. "Training" error  
e.g. "Testing" error

K

Training data can behave very differently to testing data!

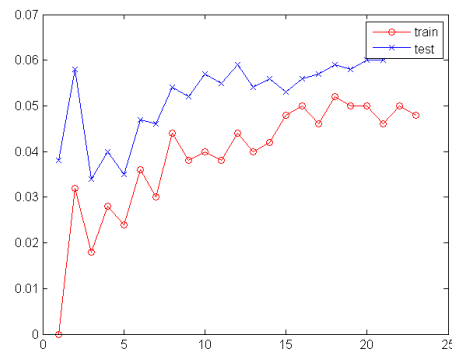
06/12/2017

EC-350 AI and DSS

EME (NUST)

41

## Classifying '3' versus '8' Digits



Best testing error at  $K=3$ , about 3.2%.

Is this "good"?

06/12/2017

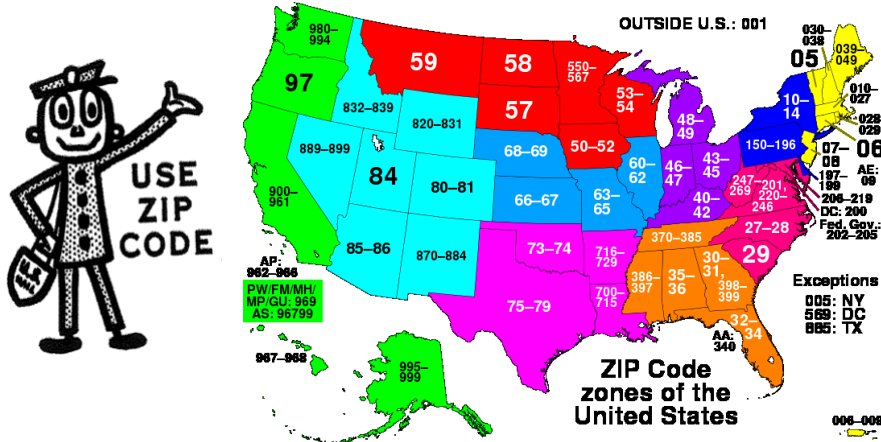
EC-350 AI and DSS

EME (NUST)

42

## Classifying '3' versus '8' Digits

Zip-codes: "8" is very common on the West Coast, while "3" is rare.  
Making a mistake will mean your Florida post ends up in Las Vegas!



06/12/2017

EC-350 AI and DSS

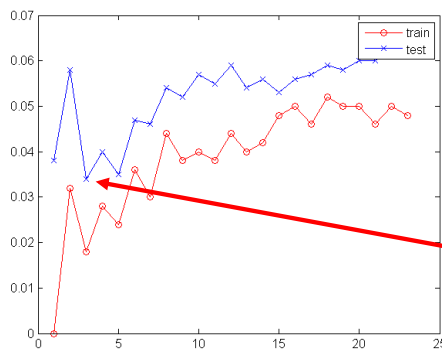
EME (NUST)

43

## Classifying '3' versus '8' Digits

Sometimes, classes are rare, so your learner will not see many of them.

What if, in testing phase you saw 1,000 digits ....



32 instances



968 instances

3.2% error was  
achieved by just saying  
"8" to everything!

06/12/2017

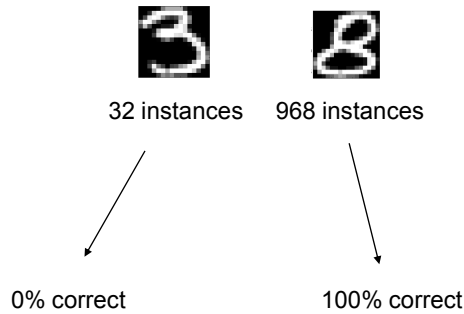
EC-350 AI and DSS

EME (NUST)

44

## Classifying '3' versus '8' Digits

### Solution?



*Measure accuracy on each class separately.*

06/12/2017

EC-350 AI and DSS

EME (NUST)

45

## R.O.C. Analysis



32 instances    968 instances

*A statistical framework.*

*Receiver Operator Characteristics  
Developed in WW-2 to assess radar operators.*

*"How good is the radar operator at spotting  
incoming bombers?"*



*False positives*

*– i.e. falsely predicting a bombing raid*

*False negatives*

*– i.e. missing an incoming bomber  
(VERY BAD!)*

06/12/2017

EC-350 AI and DSS

EME (NUST)

46

## R.O.C. Analysis

The “3” digits are like the bombers.  
Rare events but costly if we misclassify!

*False positives – i.e. falsely predicting an event*  
*False negatives – i.e. missing an incoming event*



Similarly, we have “true positives” and “true negatives”

		Prediction	
		0	1
Truth	0	TN	<b>FP</b>
	1	<b>FN</b>	TP

06/12/2017

EC-350 AI and DSS

EME (NUST)

47

## Building a “Confusion Matrix”

		Prediction	
		0	1
Truth	0	TN	<b>FP</b>
	1	<b>FN</b>	TP

$$\text{Sensitivity} = \frac{TP}{TP+FN}$$

... chances of spotting a “3” when presented with one  
(i.e. accuracy on class “3”)

$$\text{Specificity} = \frac{TN}{TN+FP}$$

... chances of spotting an 8 when presented with one  
(i.e. accuracy on class “8”)

06/12/2017

EC-350 AI and DSS

EME (NUST)

48



## R.O.C. Analysis

$$\text{Sensitivity} = \frac{TP}{TP+FN} = ?$$

$$\text{Specificity} = \frac{TN}{TN+FP} = ?$$

		Prediction	
		0	1
Truth	0	60	<b>30</b>
	1	<b>80</b>	20

TN	<b>FP</b>
<b>FN</b>	TP

60+30 = 90 examples in the dataset were class 0

80+20 = 100 examples in the dataset were class 1

90+100 = 190 examples in the data overall

06/12/2017

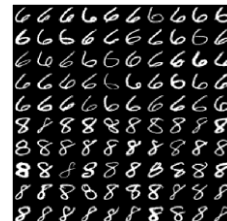
EC-350 AI and DSS

EME (NUST)

49

## Assignment # 3

- Use K-NN to recognise handwritten USPS digits.
- You will give the printed report including
  - Code
  - Results (Figures, Graphs, Tables etc.)
  - Your description (Max. 1 page)
- Deadline for submission
  - After 2 weeks



06/12/2017

EC-350 AI and DSS

EME (NUST)

50