

DEEPESH.

# K Nearest Neighbor

K, dist  
w/ me  
dataset

Jn 1. cv

distance

03 Aug 19 theory + concepts  
04 Aug 19 practical

## TIME TABLE

## ML2

1	Sat 27 July 19	PCA	
	Sun 28 July 19	PCA + practical	
2	Sat 03 Aug 19	KNN	KNN
	Sun 04 Aug 19	KNN	K means
3	Sat 10 Aug 19	K means	K means
	Sun 11 Aug 2019		
4	Sat 17 Aug 2019		
	Sun 18 Aug 2019		
5	Sat 24 Aug 2019		
	Sun 25 Aug 2019		
6	Sat 31 Aug 2019		
	Sun 01 Sep 2019		

- 2 → 1. Principal Component Analysis (PCA)
- 1 1/2 → 2. K Nearest Neighbors (KNN)
3. Naive Bayes
4. Support Vector machines (SVM)
- 1 1/2 → 5. K means (clustering)
6. Time series Analysis

≥ 20% practical

2

RECAP.PCA

principal component analysis.

variance decreasing

information

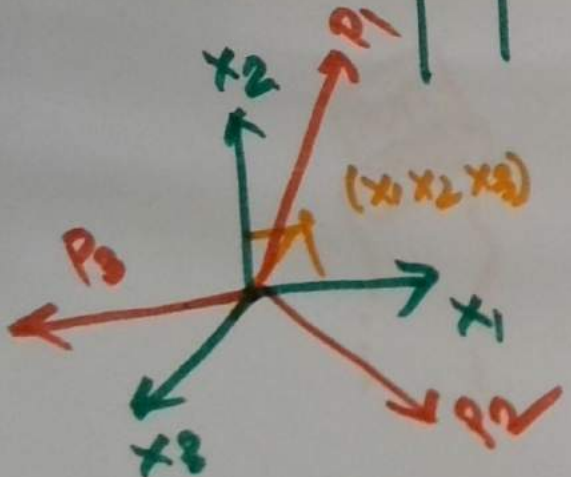
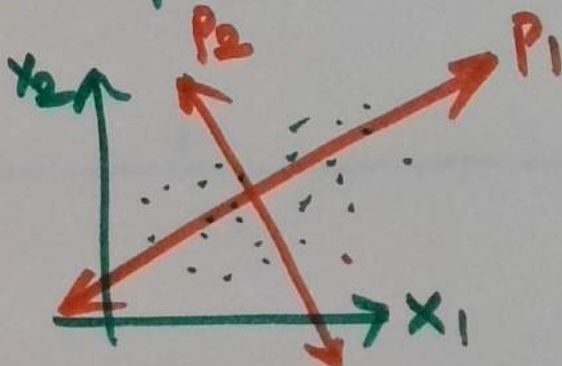
PC

$x_1$	$x_2$	$x_3$

$P_1$	$P_2$	$P_3$

drop  
last  
two

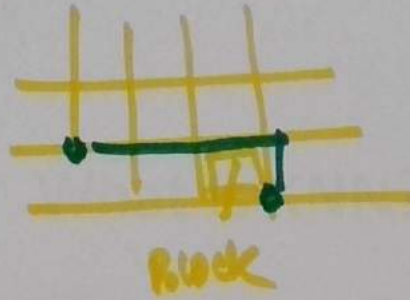
DR.



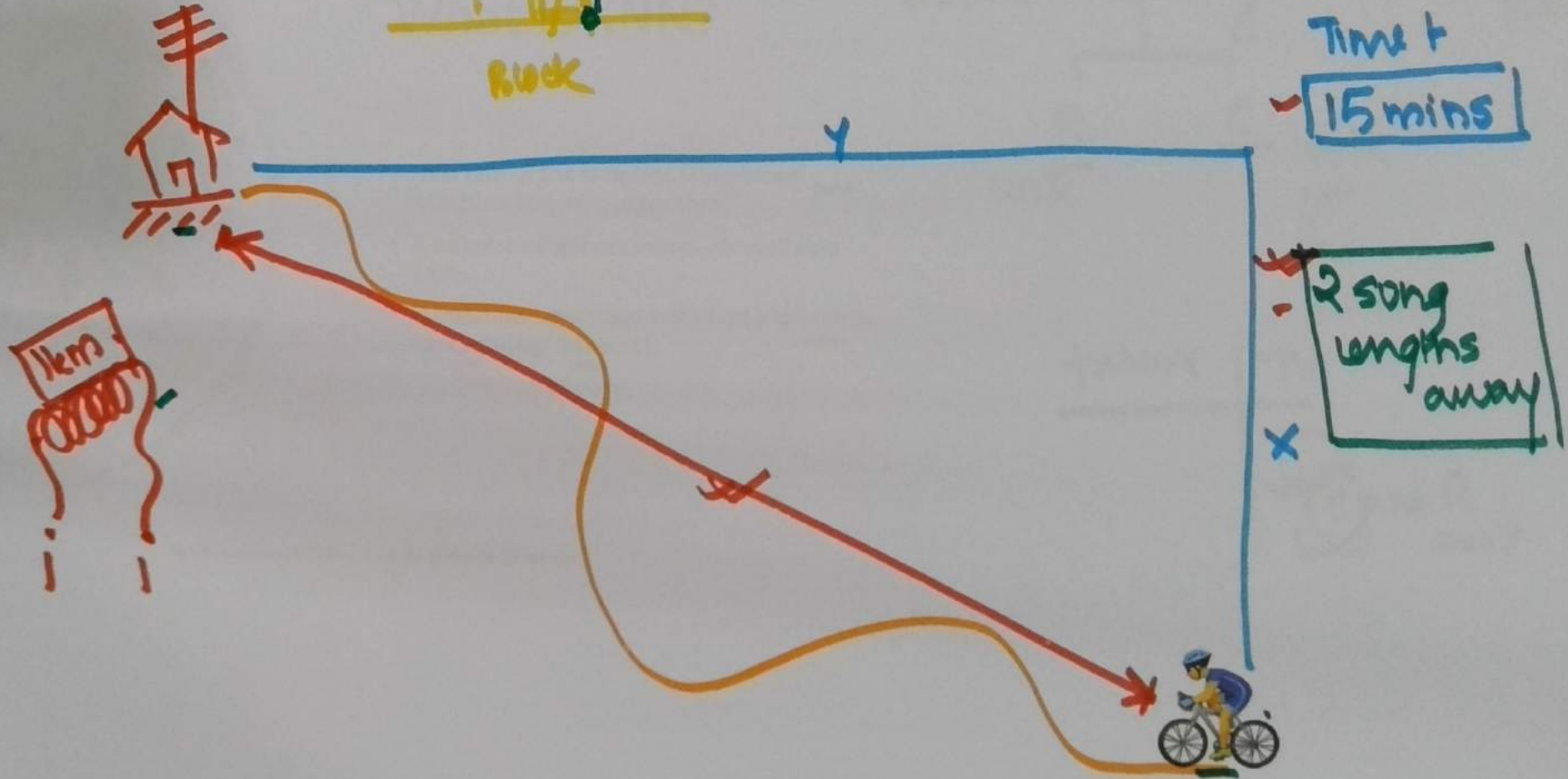


3

KNN  
└─ nearest  
└─ distance

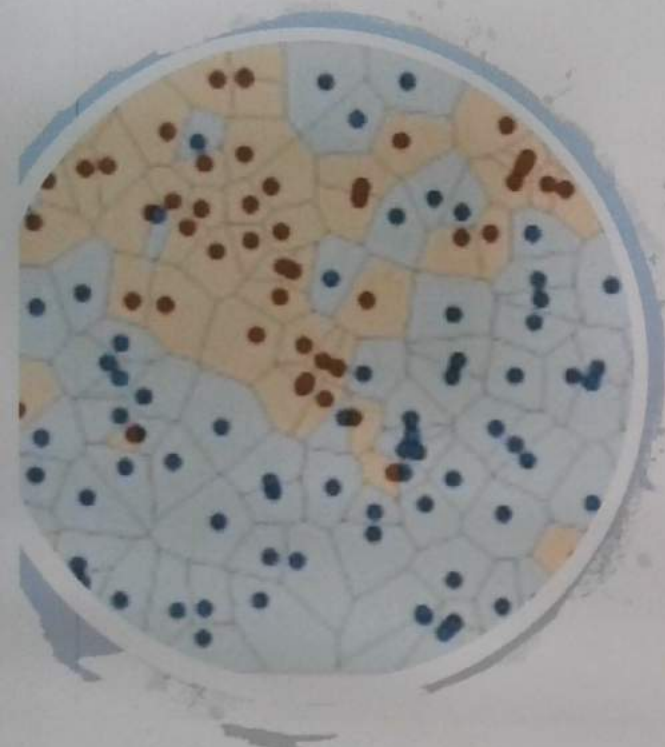


"DISTANCE"  
dissimilarity



location  
x, y  
lat/long

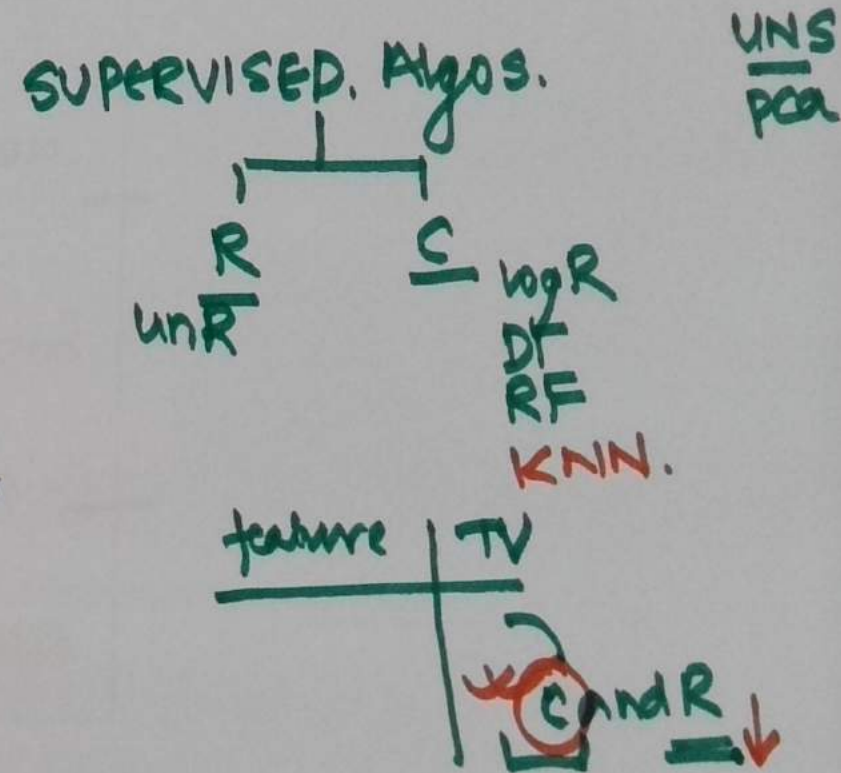
4



## What is KNN?

- It is one of the simplest **Supervised Machine Learning** algorithm.
- K nearest neighbors **stores all available cases**.
- It classifies new cases based on a **similarity measure**(eg; distance).

International School of AI & Data Science

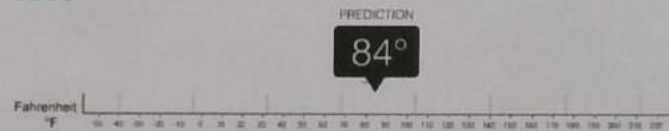


- It is used for both **classification** and **regression** predictive problems.
- However, it is more widely used in **classification problems** in the industry.



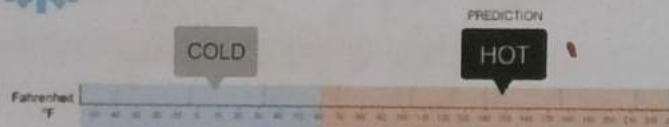
### Regression

What is the temperature going to be tomorrow?



### Classification

Will it be Cold or Hot tomorrow?



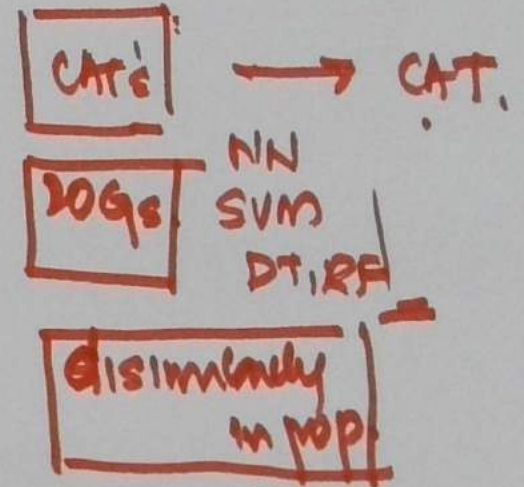


KNN is based  
on Feature  
similarity.



International School of AI & Data Science

Algorithm.



Lets consider an example:

Is that a  
dog?



No dear, you can  
differentiate  
between a cat  
and a dog based  
on their  
characteristics






length cat

sharp claws

LE	SC	DIC
8	60	D.



CATS




Sharp Claws, uses to climb

Smaller length of ears

Meows and purrs

Doesn't love to play around

DOGS

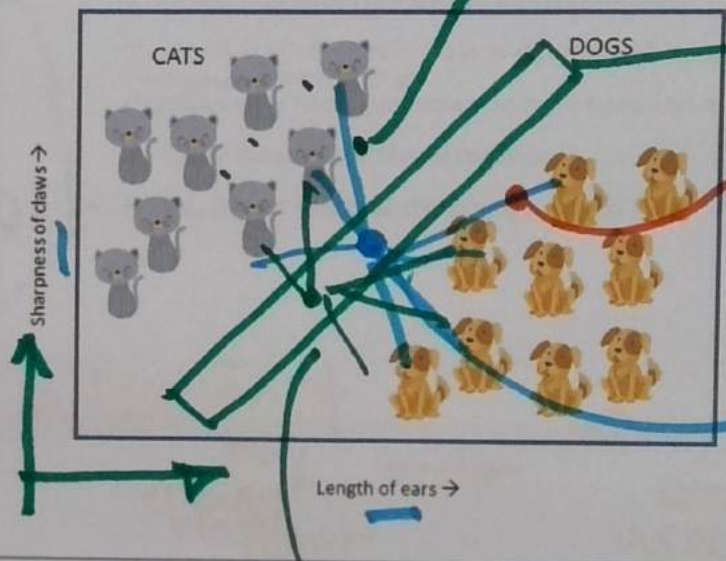


Dull Claws

Bigger length of ears

Barks

Loves to run around



unknown animal  $\equiv$  CAT

unknown  $\equiv$  DOG

unknown  $\equiv$  CAT or DOG

WHAT IF  
DISTANCE IS  
EQUAL?

check  
k of nearest  
neigh.

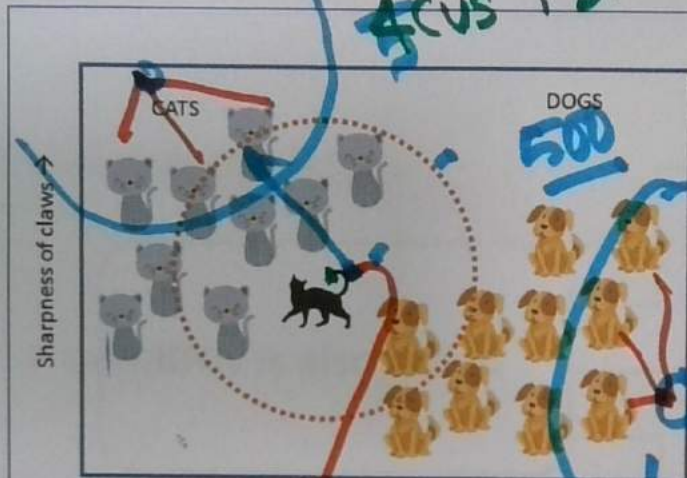
5

3/5  $\equiv$  DOG

2/5  $\equiv$  CAT

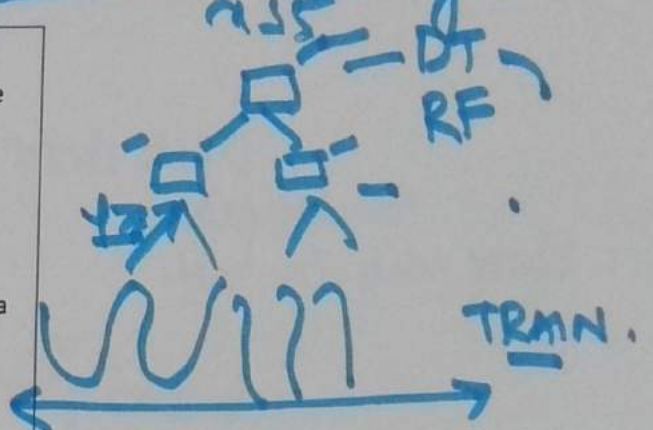
9

ACUS 1 D. = CATS



- Therefore, the features of the strange animal are more like cats.
- K nearest neighbors is a simple algorithm that stores all available cases (cats and dogs) and classifies new cases (unknown animal) based on a similarity measure or characteristics.
- Hence, it is proved that the strange animal is a cat.

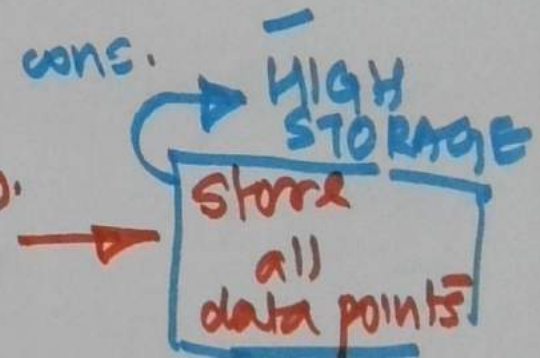
$$\text{height}(p) = \beta_0 + \beta_1 x + \beta_2 x_2 + \dots$$



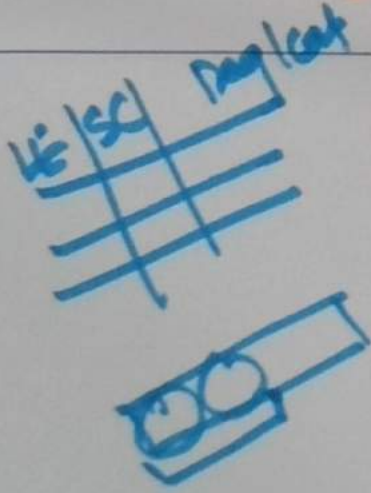
TEST

location of my TEST point

which K neigh. are used?



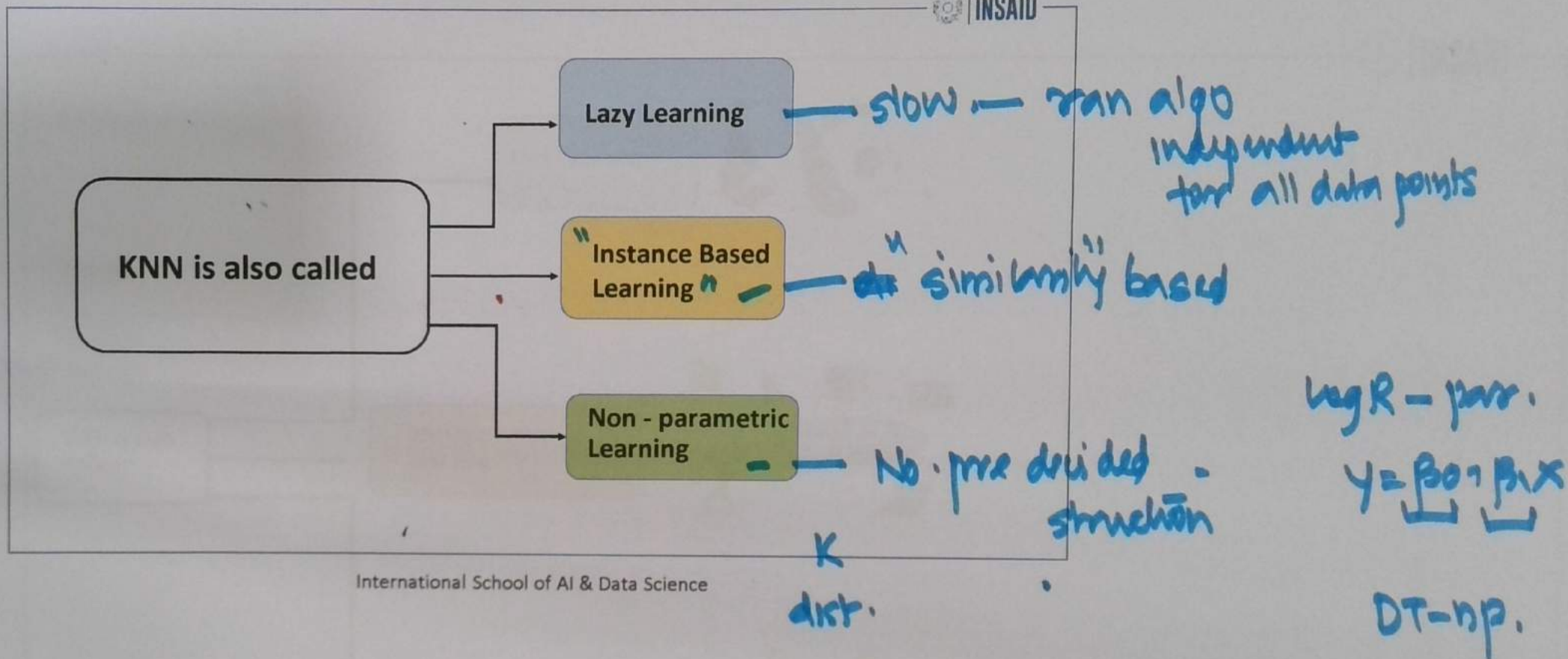
and process only after we have the TEST point



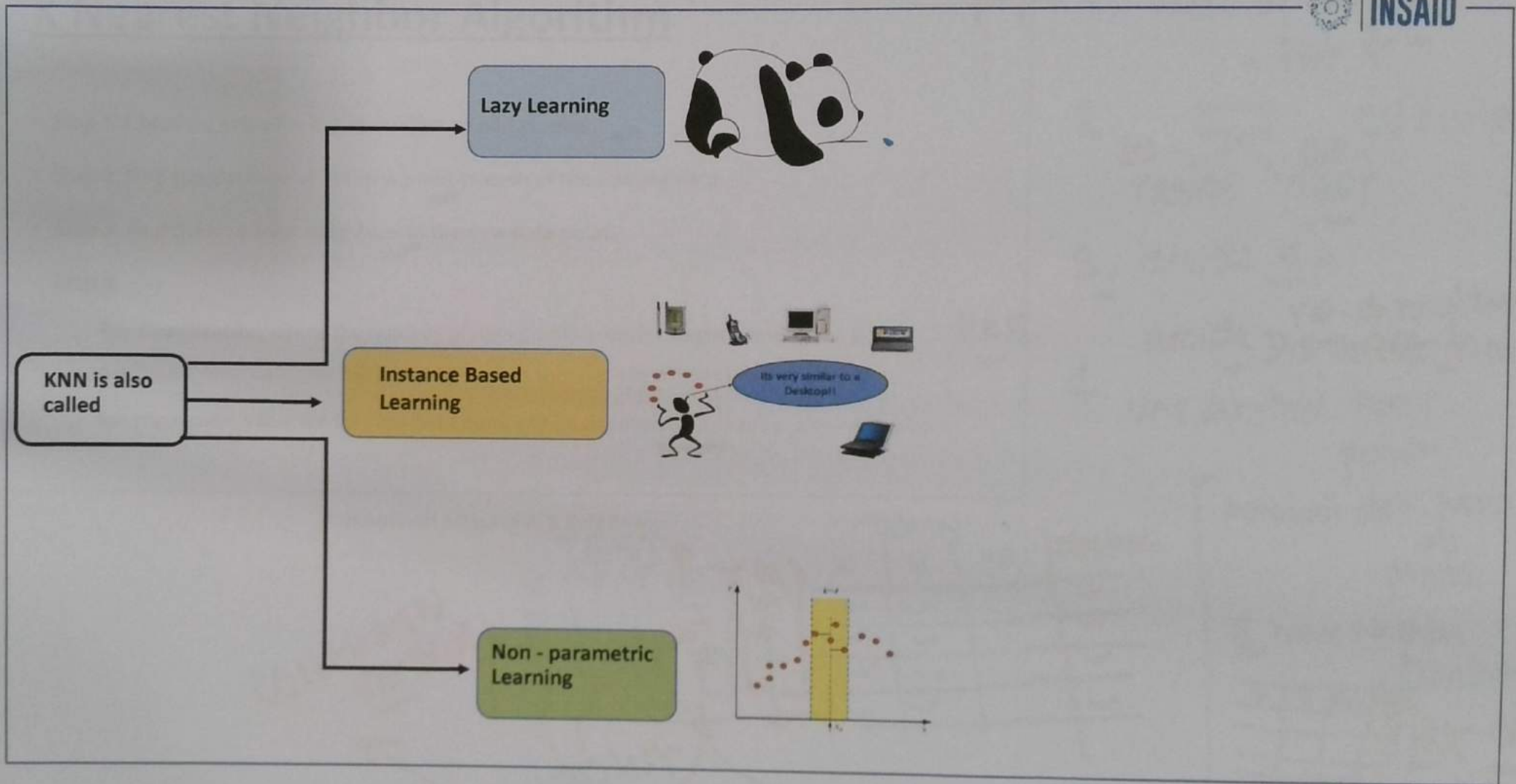
International School of AI & Data Science

6th largest dist = radius  
SLOW  
 cons










## K Nearest Neighbor Algorithm

- **Step 1:** Choose a value for K. K should be an odd number.
- **Step 2:** Find the distance of the new point to each of the training data.
- **Step 3:** Find the K nearest neighbors to the new data point.
- **Step 4:**
  - For **classification**, count the number of data points in each category among the k neighbors. New data point will belong to class that has the **most (Mode)** neighbors.
  - For **regression**, value for the new data point will be the **average (Mean)** of the k neighbors.

X	Y	cat
INSIDE		

1. Table of data -  
+ Plot it -

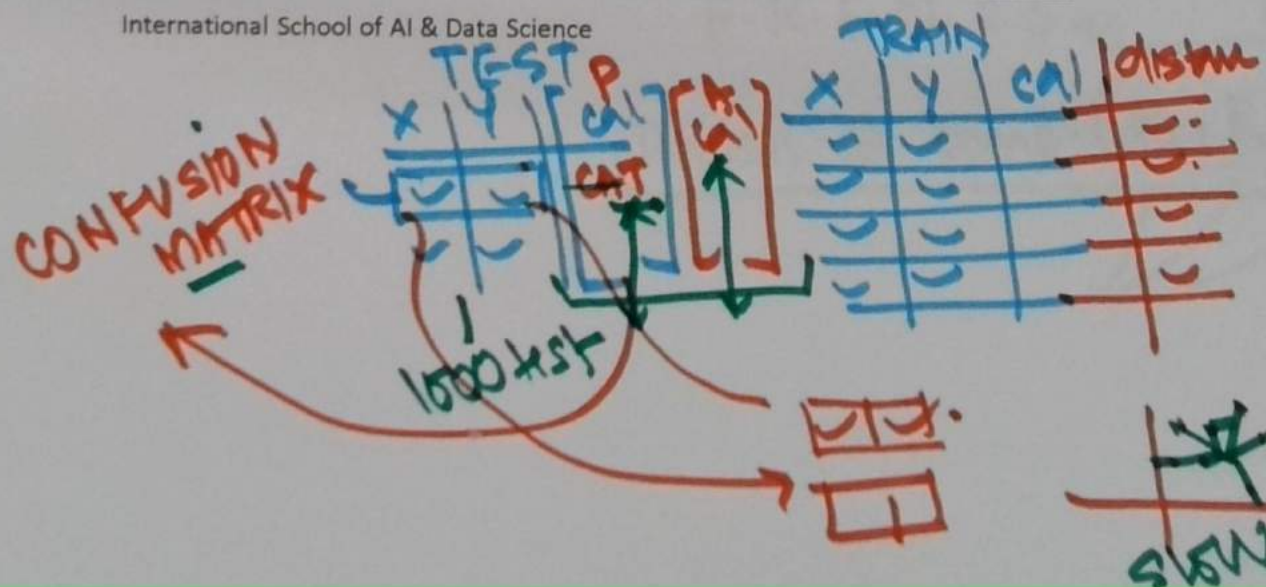
2.  CV division  
80% TRAIN 20% TEST

3. divide K -

K=5

4. one at a time TEST point  
no. of neighbours  
decide DISTANCE measure

International School of AI & Data Science



compute dist. from all points.  
5 near neighbor  
CLASSIFYING  
3/5 CAT

## What is K nearest neighbors?

- K is a number used to identify similar neighbors for the new data point.
- KNN takes K nearest neighbors to decide where the new data point will belong to.
- This decision is based on feature similarity.

### For example

- We have Friend circle in our new neighborhood.
- We select 3 neighbors that we want to be close friends based on common thinking or hobbies.
- In this case K is 3



International School of AI & Data Science

our friends

F.R.E.N.D.S

Person  $\equiv$  comb. (M, R, P, R, C, J)

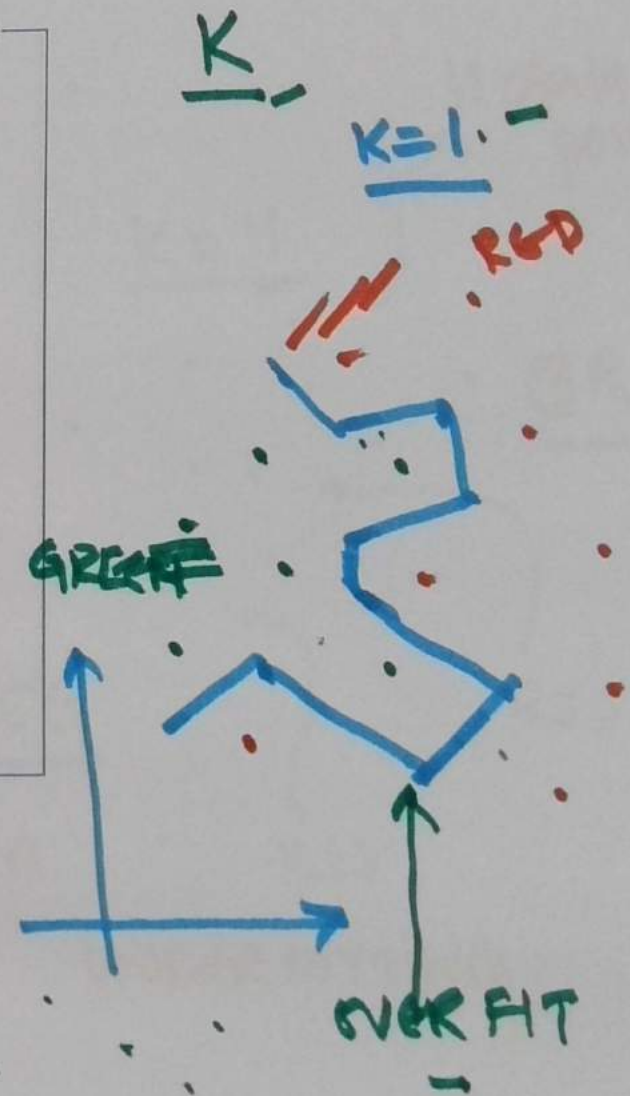


## Overfitting in KNN

- When "k" is a **very small number**- KNN can **overfit**.
- It will classify just based on the **closest neighbors** instead of learning a good separating frontier between classes.

International School of AI & Data Science

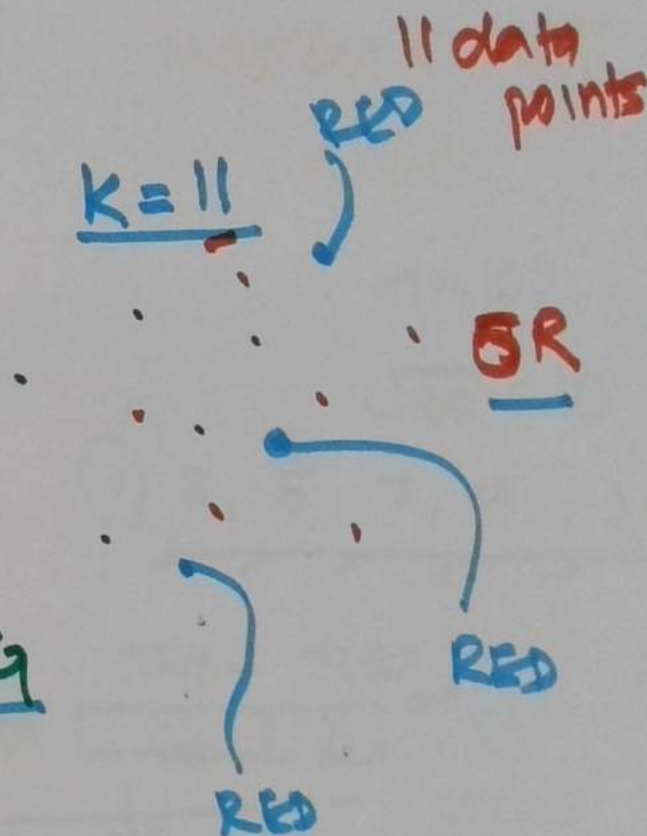
Small  $k \rightarrow$  overfitting



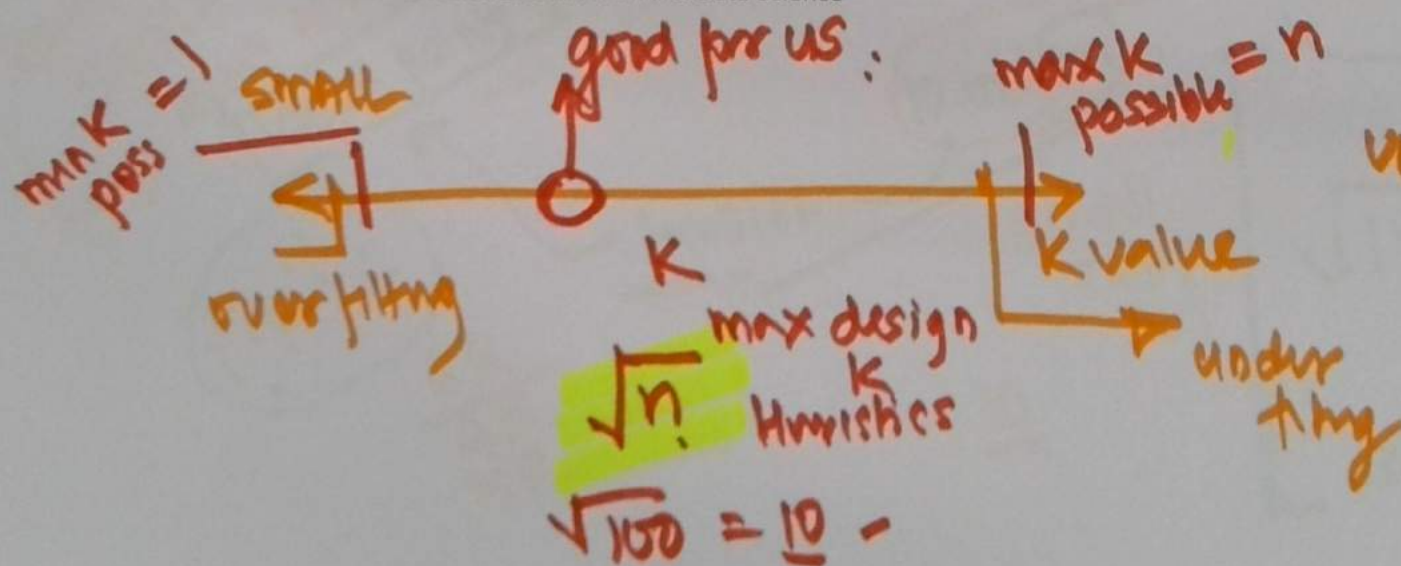
## Underfitting in KNN



- If  $k = n$ , KNN will think **every point belongs to the class that has more samples**.
- If " $k$ " is a **very big number** - KNN will **underfit**, in the limit.



International School of AI & Data Science





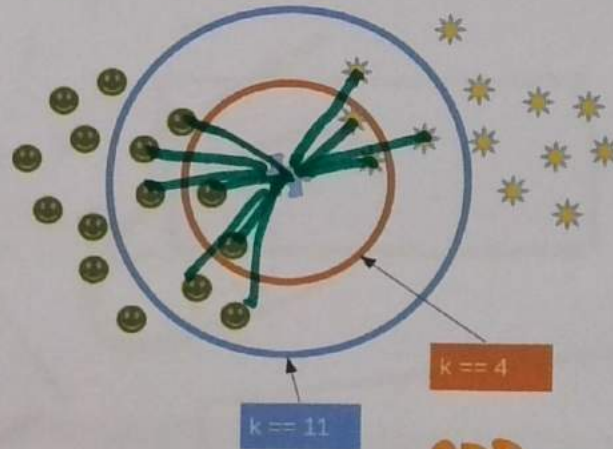
## How do we choose the factor K?

- $\text{Sqrt}(n)$ , where  $n$  is the total number of datapoints.

== or == ?

- Odd value of k is selected to avoid any confusion between two classes of data.

- In python, we will see how to use cross validation to choose the factor K (Later).



ODD

max  $K = \sqrt{n}$   
K odd.

$$n = 100$$

$$\sqrt{100} = 10$$

① 3, 5, 7, 9, 11

TRAIN TEST

80 20

CV  
cross  
validation

$$\sqrt{\frac{(x_2 - x_1)^2 + (y_2 - y_1)^2}{(x_2 - x_1)^2 + \dots}}$$

confusion. ==

EASY to EXPLAIN.  
SIMPLE  
DEANT ACCURACY.  
10000 points  
6R  
4B  
5R

International School of AI & Data Science

1 million

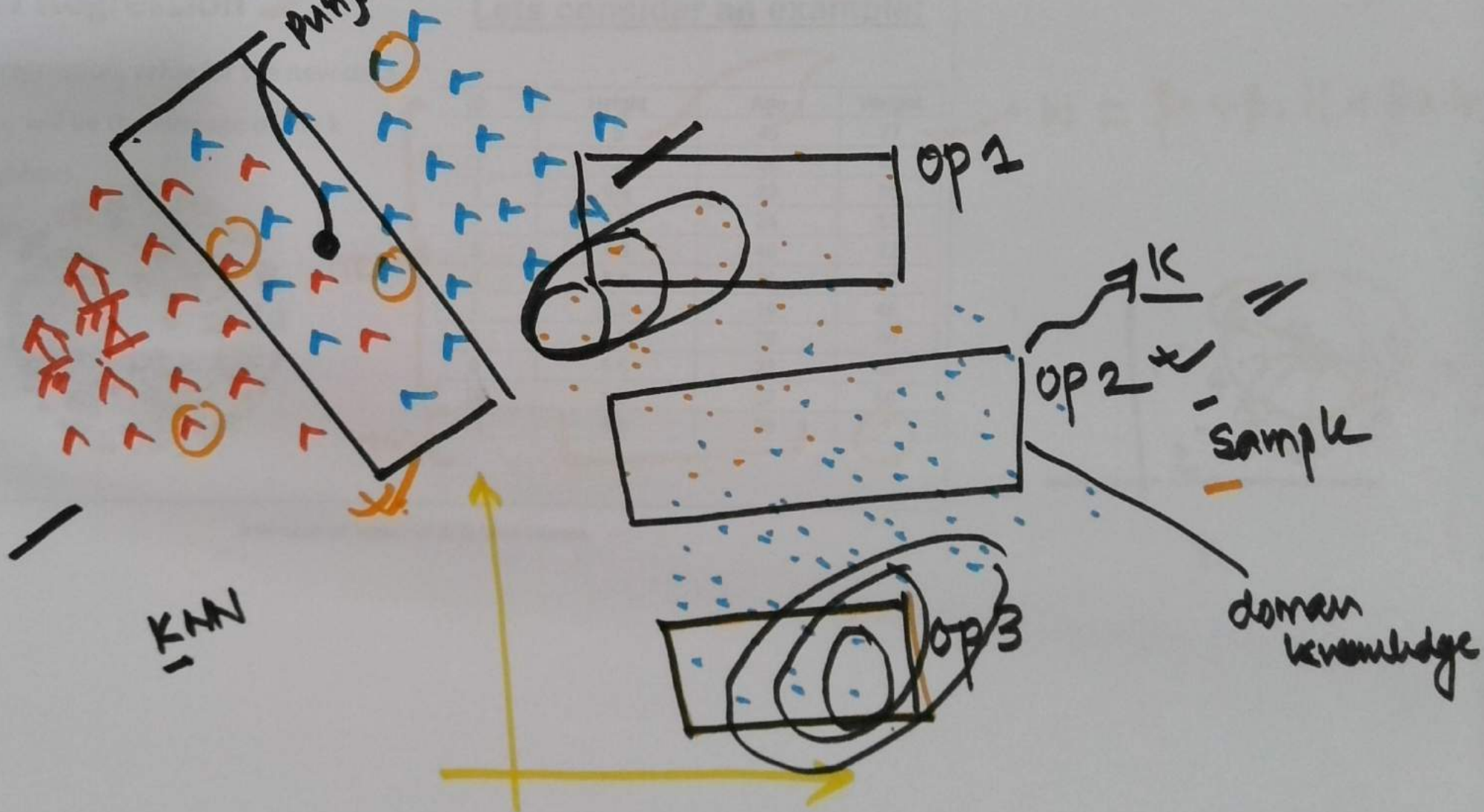
10 million

Small TEST SET

10



पुण्याची वा  
सुगली?



## KNN Regression

- For regression, value for the **new data point** will be **the average of the k neighbors**.



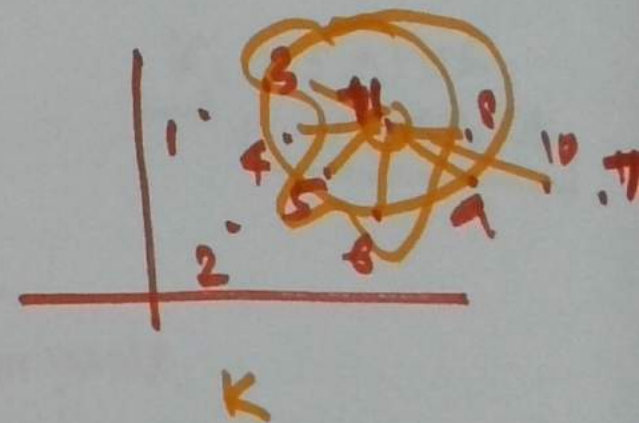
TRAIN

TEST

Lets consider an example:

ID	Height	Age	Weight
1	5	45	77
2	5.11	26	47
3	5.6	30	55
4	5.9	34	59
5	4.8	40	72
6	5.8	36	60
7	5.3	19	40
8	5.8	28	60
9	5.5	23	45
10	5.6	32	58
11	5.5	38	?

$$W = \beta_0 + \beta_1 H + \beta_2 A$$





19

K=5

5 45  
5.5, 38

INSAD

ID	Height	Age	Weight	Distance
1	5	45	77	7.018
2	5.11	26	47	
3	5.6	30	55	
4	5.9	34	59	
5	4.8	40	72	
6	5.8	36	60	
7	5.3	19	40	
8	5.8	28	60	
9	5.5	23	45	
10	5.6	32	58	
11	5.5	38	?	

distance

$$\sqrt{(5.5 - 5)^2 + (45 - 38)^2} = 7.018$$

K, distance measure

5 nearest data points

5NN

International School of AI &amp; Data Science

Regression

52.8

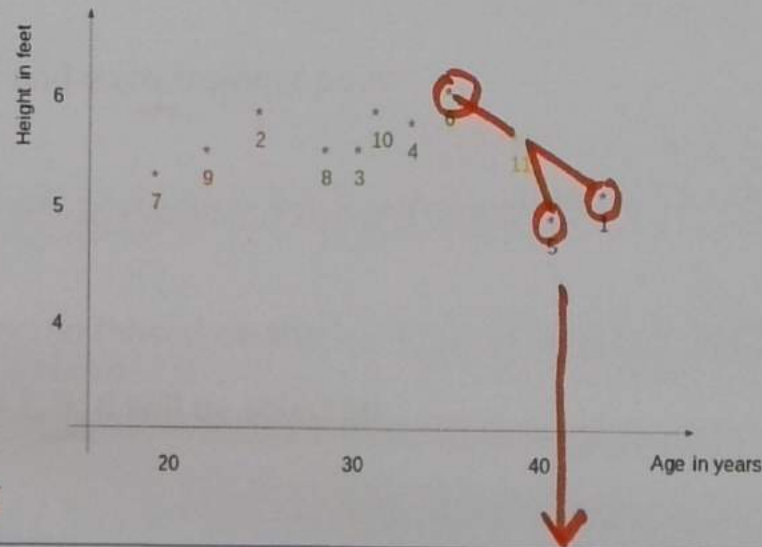
TEST 12  
13  
14

TRAIN



Below is the plot of height versus age from the table:

- Y-axis represents the **height of a person** (in feet)
- X-axis represents the **age** (in years).
- The points are numbered according to the **ID values**.
- The yellow point (ID 11) is our **test point**.



International School of AI & Data Science

$$\frac{1 \ 2 \ 3 \ 4}{\hline} \\ \approx 3.57$$

average  $\equiv$  REG.

calculate all 10 distances

↓  
shortest 3

## Step 1

Distance between the new point and each training point is calculated.

## Step 2

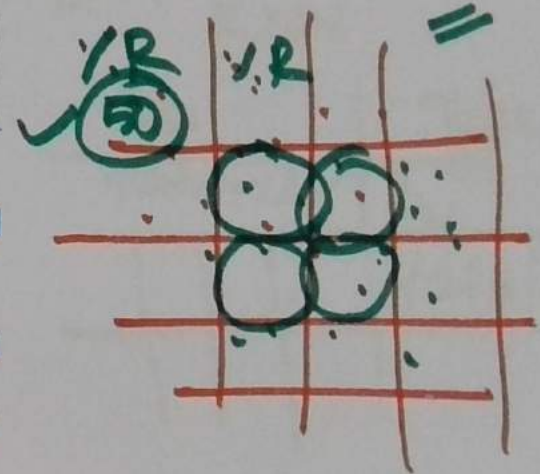
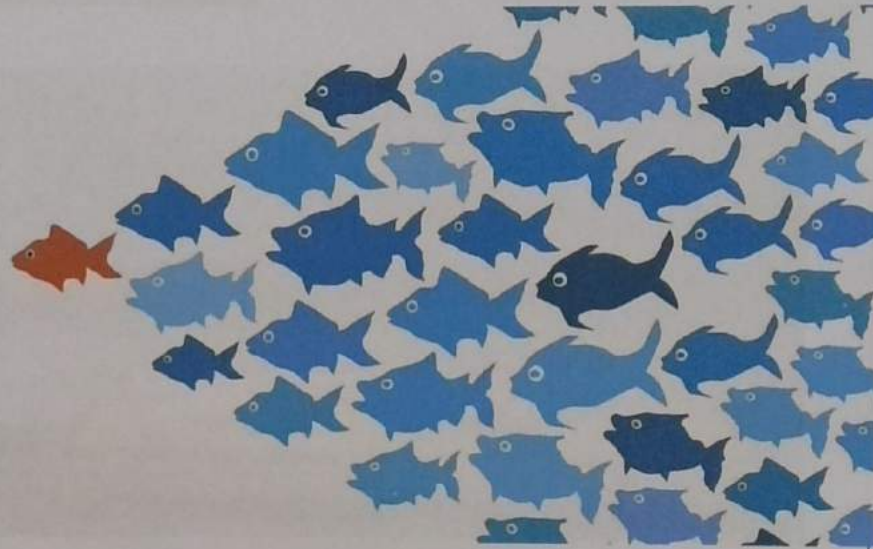
The closest k data points are selected (based on the distance). In this example, points 1, 5, 6 will be selected if value of k is 3

## Step 3

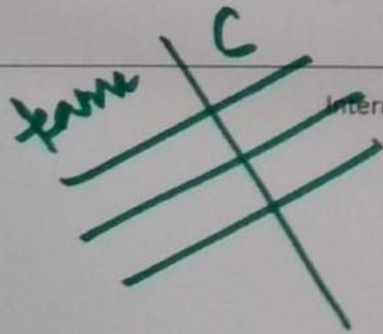
The average of these data points is the final prediction for the new point. Here, we have weight of ID11 =  $(77+72+60)/3 = 69.66$  kg.

# KNN Classification

- For classification, count the number of data points in each category among the k neighbors.
- New data point will belong to class that has the most neighbors.



5/1 X

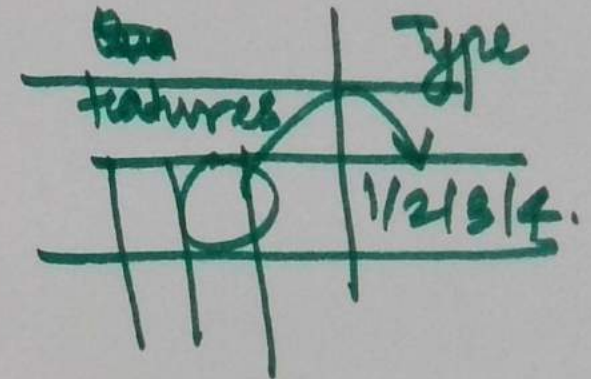




Wine Quality database -  
KNN

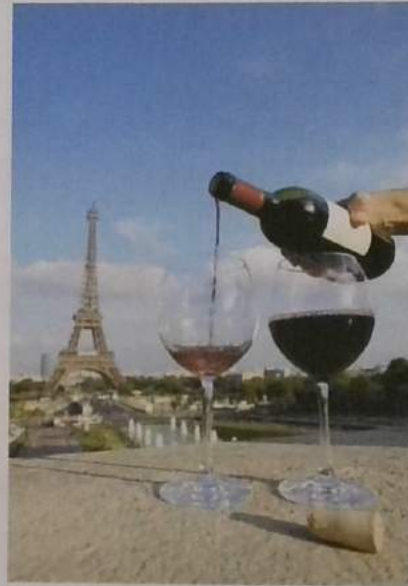
Lets consider an example:

Wine Quality  
Detection with KNN  
Algorithm



## What if want to classify a new wine quality ?

- In this case , we would find the distance between this wine feature with all the wines.



### Wine Dataset

	fixed acidity	volatile acidity	citric acid	residual sugar	chlorides	free sulfur dioxide	total sulfur dioxide	density	pH	sulphates	alcohol	quality
0	7.4	0.70	0.00	1.9	0.076	11.0	34.0	0.9978	3.51	0.56	9.4	5
1	7.8	0.88	0.00	2.6	0.098	25.0	67.0	0.9968	3.20	0.68	9.8	5
2	7.8	0.76	0.04	2.3	0.092	15.0	54.0	0.9970	3.26	0.65	9.8	5
3	11.2	0.28	0.56	1.9	0.075	17.0	60.0	0.9980	3.16	0.58	9.8	6
4	7.4	0.70	0.00	1.9	0.076	11.0	34.0	0.9978	3.51	0.56	9.4	5

- The dataset comprise of 12 features.
- We classified wine on the basis of closeness of different features.



TV

Experiment K =

Ardur

K

Distances =



## Euclidean distance

is calculated as the square root of the sum of the squared differences between a new point (x) and an existing point (y).

$$\sqrt{\sum_{i=1}^k (x_i - y_i)^2}$$

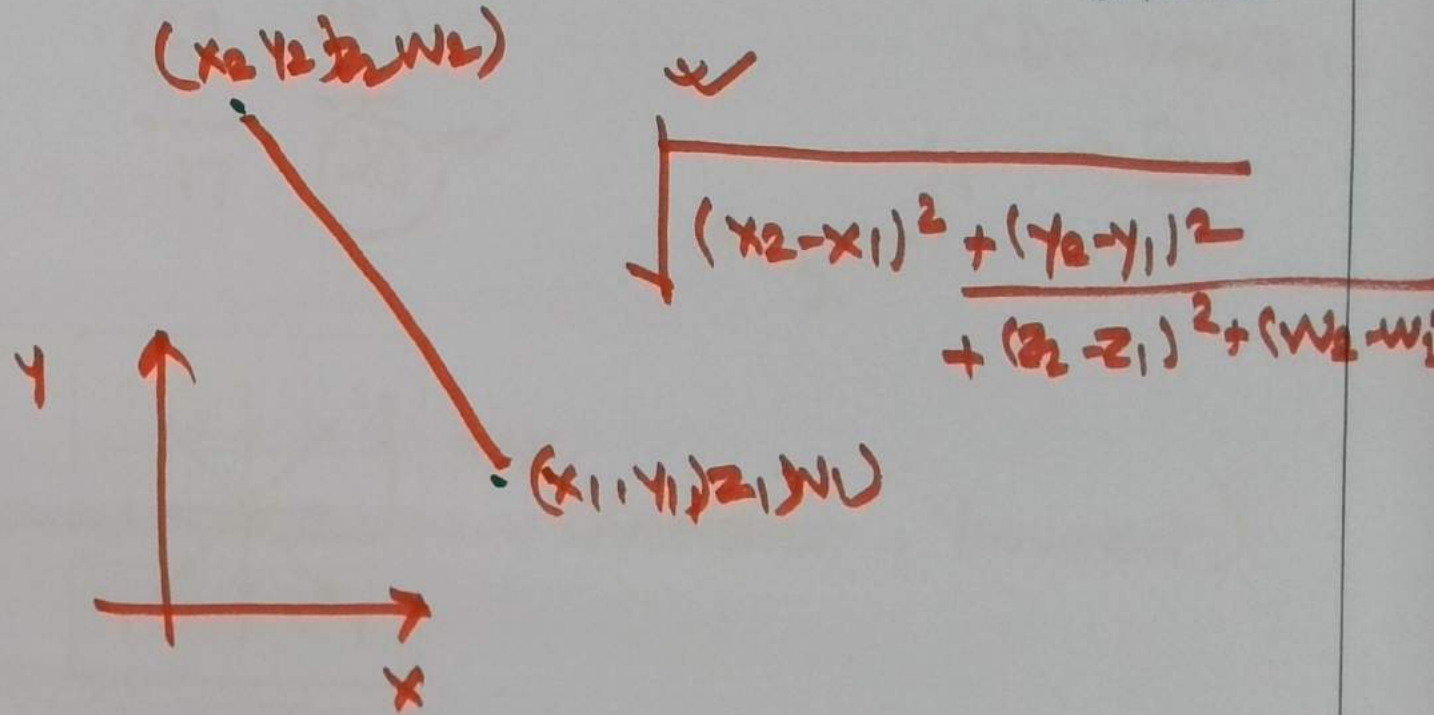
~~$\sqrt{\sum_{i=1}^k (x_i - y_i)^2}$~~

~~$\sqrt{\sum_{i=1}^k (x_i - y_i)^2}$~~

~~$\sqrt{\sum_{i=1}^k (x_i - y_i)^2}$~~

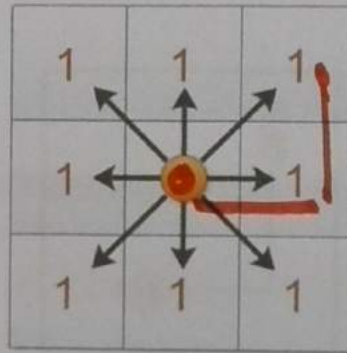
~~$\sqrt{\sum_{i=1}^k (x_i - y_i)^2}$~~

~~$\sqrt{\sum_{i=1}^k (x_i - y_i)^2}$~~



## Chebyshev distance

maximum metric, or  $L_\infty$  metric is a metric defined on a vector space where the **distance** between two vectors is the greatest of their differences along any coordinate dimension.



$$\max(|\underline{x_1} - \underline{x_2}|, |\underline{y_1} - \underline{y_2}|)$$

$(45, 5)$

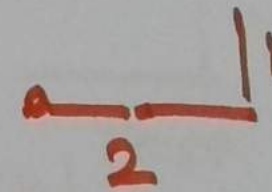
$(62, 25)$

$\frac{17}{20}$  ✓

INSAD

$CD = \max(2, 1)$

$= 2$



$\max(x\text{movement}, y\text{movement})$

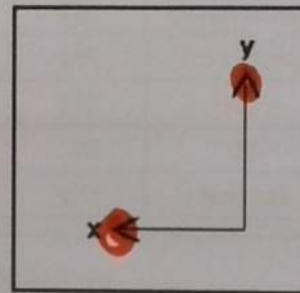


## Manhattan distance

between two vectors (or points) a and b is defined as

$$D = \sum_{i=1}^n |x_i - y_i|$$

This is known as Manhattan distance because all paths from the bottom left to top right of this idealized city have the same distance:



Manhattan

$$(x_1, y_1) \quad (x_2, y_2)$$

$$|x_1 - x_2| + |y_1 - y_2|$$

$x_{mov}$   $y_{mov}$

$$x_{mov} + y_{mov}$$

$$(45, 5) \quad (62, 25)$$

17 20

$$MD = 17 + 20 = \underline{37}$$



## When do we use KNN?

We can use K Nearest Neighbor when

Dataset is small



Because KNN is a "lazy learner"

Data is labelled



Dog

Dataset is noise free

Weight	Height	Class
51	182	Underweight
62	165	One - forty
69	176	Hello
64	173	23
65	172	Normal

Noise

CONV  
TV  
SUPERVISED.

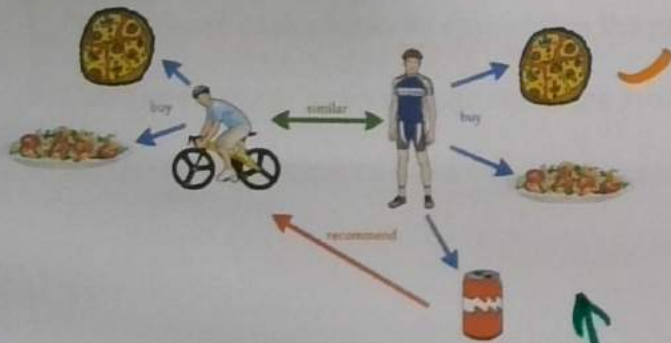
#  
NOISE  
FREE

SMALL  
DATASET

## Applications

### Recommender system

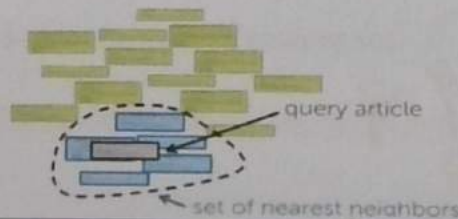
If you know a user likes a particular item, then you can **recommend similar items** for them.



### Document Retrieval

We compute **distances from query article** to all other articles. Then we search for the articles with **smallest distance** to the query article. They are called **nearest neighbors**.

Space of all articles,  
organized by similarity of text



International School of AI & Data Science

AIUP movies  
fat  
CSV  
XLS

Thriller  
Thriller

Amo  
Kw  
Bm



## Advantages of K-nearest neighbors algorithm

- It is simple to implement. \*
- It executes quickly for small training data sets.
- Performance asymptotically approaches the performance of the Bayes Classifier.
- Don't need any prior knowledge about the structure of data in the training set.
- No retraining is required if the new training pattern is added to the existing training set.

Naive Bayes

← BAYES CLASSIFICATION

$P_0, P_1$   
↙ ↘

## Limitations of K-nearest neighbors algorithm



- When the training set is large, it may take a lot of space.
- For every test data, the distance should be computed between test data and all the training data. Thus a lot of time may be needed for the testing.