# Chapter 4 :Supervised Learning

# Supervised vs. Unsupervised Learning

- Supervised learning (classification)

  - Supervision: The training data (observations, measurements, etc.) are accompanied by **labels** indicating the class of the observations

  - New data is classified based on the training set

- Unsupervised learning (clustering)

  - The class labels of training data is unknown

  - Given a set of measurements, observations, etc. with the aim of establishing the existence of classes or clusters in the data

# Prediction Problems: Classification vs. Numeric Prediction

- Classification
  - predicts categorical class labels (discrete or nominal)
  - classifies data (constructs a model) based on the training set and the values (class labels) in a classifying attribute and uses it in classifying new data
- Numeric Prediction
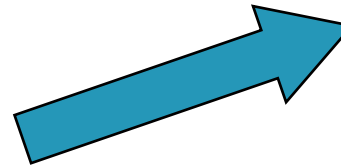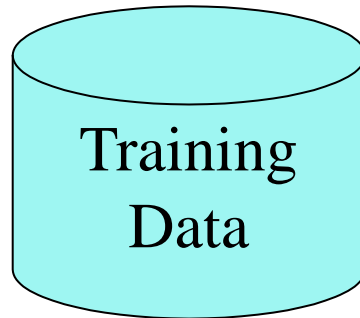  - models continuous-valued functions, i.e., predicts unknown or missing values
- Typical applications
  - Credit/loan approval:
  - Medical diagnosis: if a tumor is cancerous or benign
  - Fraud detection: if a transaction is fraudulent
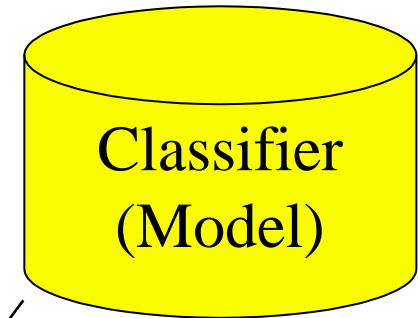  - Web page categorization: which category it is

# Classification—A Two-Step Process

- Model construction: describing a set of predetermined classes
  - Each tuple/sample is assumed to belong to a predefined class, as determined by the class label attribute
  - The set of tuples used for model construction is training set
  - The model is represented as classification rules, decision trees, or mathematical formulae
- Model usage: for classifying future or unknown objects
  - Estimate accuracy of the model
    - The known label of test sample is compared with the classified result from the model
    - Accuracy rate is the percentage of test set samples that are correctly classified by the model
    - Test set is independent of training set (otherwise overfitting)
  - If the accuracy is acceptable, use the model to classify new data
- Note: If *the test set* is used to select models, it is called validation (test) set

# Process (1): Model Construction



Training Data

Classification Algorithms

Classifier (Model)

| NAME | RANK | YEARS | TENURED |
|------|------|-------|---------|
| Mike | Assistant Prof | 3 | no |
| Mary | Assistant Prof | 7 | yes |
| Bill | Professor | 2 | yes |
| Jim | Associate Prof | 7 | yes |
| Dave | Assistant Prof | 6 | no |
| Anne | Associate Prof | 3 | no |

IF rank = 'professor'
OR years > 6
THEN tenured = 'yes'

# Process (2): Using the Model in Prediction



Classifier

Testing Data

Unseen Data

(Jeff, Professor, 4)

| NAME | RANK | YEARS | TENURED |
|------|------|-------|---------|
| Tom | Assistant Prof | 2 | no |
| Merlisa | Associate Prof | 7 | no |
| George | Professor | 5 | yes |
| Joseph | Assistant Prof | 7 | yes |

Tenured?

Yes

# Naïve Bayes

# Bayes' Theorem: Basics

- Total probability Theorem: $P(B) = \sum_{i=1}^{M} P(B|A_i)P(A_i)$

- Bayes' Theorem: $$P(H|\mathbf{X}) = \frac{P(\mathbf{X}|H)P(H)}{P(\mathbf{X})} = P(\mathbf{X}|H) \times P(H) / P(\mathbf{X})$$

  - Let $\mathbf{X}$ be a data sample ("*evidence*"): class label is unknown
  - Let H be a *hypothesis* that X belongs to class C
  - Classification is to determine $P(H|\mathbf{X})$, (i.e., *posteriori probability):* the probability that the hypothesis holds given the observed data sample $\mathbf{X}$
  - P(H) (*prior probability*): the initial probability
    - E.g., $\mathbf{X}$ will buy computer, regardless of age, income, …
  - P($\mathbf{X}$): probability that sample data is observed
  - P($\mathbf{X}$|H) (likelihood): the probability of observing the sample $\mathbf{X}$, given that the hypothesis holds
    - E.g., Given that $\mathbf{X}$ will buy computer, the prob. that X is 31..40, medium income

# Prediction Based on Bayes' Theorem

- Given training data **X**, *posteriori probability of a hypothesis* H, P(H|**X**), follows the Bayes' theorem

$$P(H \mid \mathbf{X}) = \frac{P(\mathbf{X} \mid H)P(H)}{P(\mathbf{X})} = P(\mathbf{X} \mid H) \times P(H) / P(\mathbf{X})$$

- Informally, this can be viewed as

    posteriori = likelihood x prior/evidence

- Predicts **X** belongs to $C_i$ iff the probability $P(C_i \mid \mathbf{X})$ is the highest among all the $P(C_k \mid X)$ for all the *k* classes

- Practical difficulty: It requires initial knowledge of many probabilities, involving significant computational cost

# Classification Is to Derive the Maximum Posteriori

- Let D be a training set of tuples and their associated class labels, and each tuple is represented by an n-D attribute vector $\mathbf{X} = (x_1, x_2, \ldots, x_n)$

- Suppose there are $m$ classes $C_1, C_2, \ldots, C_m$.

- Classification is to derive the maximum posteriori, i.e., the maximal $P(C_i | \mathbf{X})$

- This can be derived from Bayes' theorem

$$P(C_i | \mathbf{X}) = \frac{P(\mathbf{X} | C_i) P(C_i)}{P(\mathbf{X})}$$

- Since P(X) is constant for all classes, only

  needs to be maximized
$$P(C_i | \mathbf{X}) = P(\mathbf{X} | C_i) P(C_i)$$

# Naïve Bayes Classifier

- A simplified assumption: attributes are conditionally independent (i.e., no dependence relation between attributes):

$$P(\mathbf{X} \mid C_i) = \prod_{k=1}^{n} P(x_k \mid C_i) = P(x_1 \mid C_i) \times P(x_2 \mid C_i) \times \ldots \times P(x_n \mid C_i)$$

- This greatly reduces the computation cost: Only counts the class distribution

- If $A_k$ is categorical, $P(x_k \mid C_i)$ is the # of tuples in $C_i$ having value $x_k$ for $A_k$ divided by $|C_{i,D}|$ (# of tuples of $C_i$ in D)

- If $A_k$ is continous-valued, $P(x_k \mid C_i)$ is usually computed based on Gaussian distribution with a mean $\mu$ and standard deviation $\sigma$

$$g(x, \mu, \sigma) = \frac{1}{\sqrt{2\pi}\sigma} e^{-\frac{(x-\mu)^2}{2\sigma^2}}$$

and $P(x_k \mid C_i)$ is

$$P(\mathbf{X} \mid C_i) = g(x_k, \mu_{C_i}, \sigma_{C_i})$$

# Naïve Bayes Classifier: Training Dataset

Class:

C1:buys_computer = 'yes'

C2:buys_computer = 'no'

Data to be classified:

X = (age <=30,

Income = medium,

Student = yes

Credit_rating = Fair)

| age | income | student | credit_rating | com |
|------|--------|---------|---------------|-----|
| <=30 | high | no | fair | no |
| <=30 | high | no | excellent | no |
| 31…40 | high | no | fair | yes |
| >40 | medium | no | fair | yes |
| >40 | low | yes | fair | yes |
| >40 | low | yes | excellent | no |
| 31…40 | low | yes | excellent | yes |
| <=30 | medium | no | fair | no |
| <=30 | low | yes | fair | yes |
| >40 | medium | yes | fair | yes |
| <=30 | medium | yes | excellent | yes |
| 31…40 | medium | no | excellent | yes |
| 31…40 | high | yes | fair | yes |
| >40 | medium | no | excellent | no |

# Naïve Bayes Classifier: An Example

| age | income | student | credit_rating | com |
|------|--------|---------|---------------|-----|
| <=30 | high | no | fair | no |
| <=30 | high | no | excellent | no |
| 31…40 | high | no | fair | yes |
| >40 | medium | no | fair | yes |
| >40 | low | yes | fair | yes |
| >40 | low | yes | excellent | no |
| 31…40 | low | yes | excellent | yes |
| <=30 | medium | no | fair | no |
| <=30 | low | yes | fair | yes |
| >40 | medium | yes | fair | yes |
| <=30 | medium | yes | excellent | yes |
| 31…40 | medium | no | excellent | yes |
| 31…40 | high | yes | fair | yes |
| >40 | medium | no | excellent | no |

- $P(C_i)$:    $P(\text{buys\_computer} = \text{"yes"}) = 9/14 = 0.643$
  $P(\text{buys\_computer} = \text{"no"}) = 5/14 = 0.357$

- Compute $P(X|C_i)$ for each class
  $P(\text{age} = \text{"<=30"} \mid \text{buys\_computer} = \text{"yes"}) = 2/9 = 0.222$
  $P(\text{age} = \text{"<= 30"} \mid \text{buys\_computer} = \text{"no"}) = 3/5 = 0.6$
  $P(\text{income} = \text{"medium"} \mid \text{buys\_computer} = \text{"yes"}) = 4/9 = 0.444$
  $P(\text{income} = \text{"medium"} \mid \text{buys\_computer} = \text{"no"}) = 2/5 = 0.4$
  $P(\text{student} = \text{"yes"} \mid \text{buys\_computer} = \text{"yes"}) = 6/9 = 0.667$
  $P(\text{student} = \text{"yes"} \mid \text{buys\_computer} = \text{"no"}) = 1/5 = 0.2$
  $P(\text{credit\_rating} = \text{"fair"} \mid \text{buys\_computer} = \text{"yes"}) = 6/9 = 0.667$
  $P(\text{credit\_rating} = \text{"fair"} \mid \text{buys\_computer} = \text{"no"}) = 2/5 = 0.4$

- **X = (age <= 30 , income = medium, student = yes, credit_rating = fair)**

**$P(X|C_i)$ :** $P(X|\text{buys\_computer} = \text{"yes"}) = 0.222 \times 0.444 \times 0.667 \times 0.667 = 0.044$
  $P(X|\text{buys\_computer} = \text{"no"}) = 0.6 \times 0.4 \times 0.2 \times 0.4 = 0.019$

**$P(X|C_i)*P(C_i)$ :** $P(X|\text{buys\_computer} = \text{"yes"}) * P(\text{buys\_computer} = \text{"yes"}) = 0.028$
  $P(X|\text{buys\_computer} = \text{"no"}) * P(\text{buys\_computer} = \text{"no"}) = 0.007$

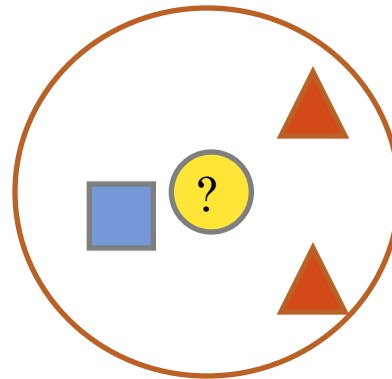**Therefore,  X belongs to class ("buys_computer = yes")**

# KNN

# WHY NEAREST NEIGHBOR?

- Used to classify objects based on closest training examples in the feature space
  - Feature space: raw data transformed into sample vectors of fixed length using feature extraction (Training Data)
- Top 10 Data Mining Algorithm
  - ICDM paper – December 2007
- Among the simplest of all Data Mining Algorithms
  - Classification Method
- Implementation of lazy learner
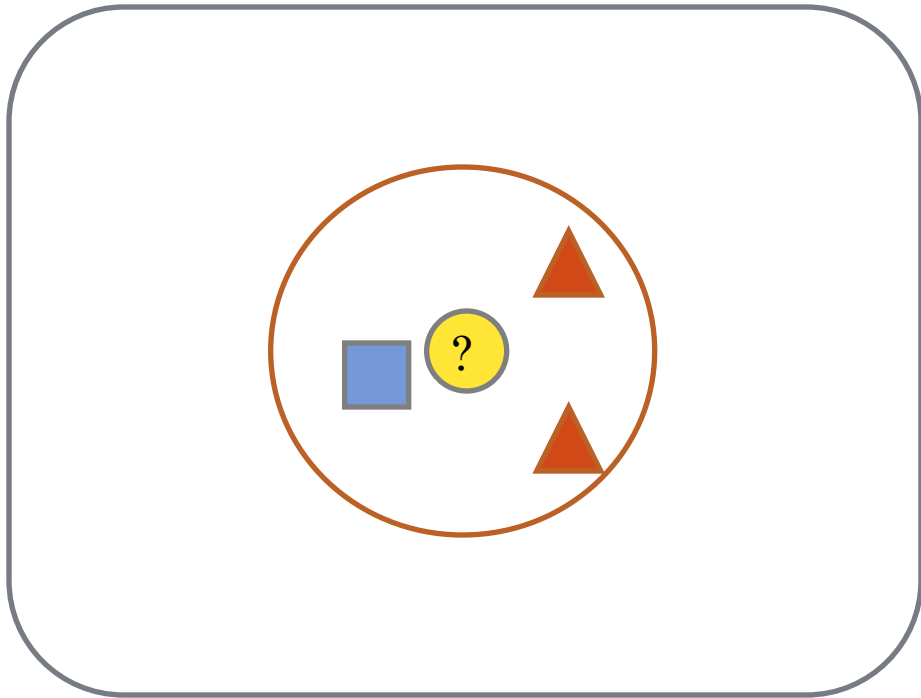  - All computation deferred until                classification

# NEAREST NEIGHBOR CLASSIFICATION

- Nearest Neighbor Overview

- *k* Nearest Neighbor

# *k* NEAREST NEIGHBOR



- Requires 3 things:
  - Feature Space(Training Data)
  - Distance metric
    - to compute distance between records
  - The value of *k*
    - the number of nearest neighbors to retrieve from which to get majority class
- To classify an unknown record:
  - Compute distance to other training records
  - Identify *k* nearest neighbors
  - Use class labels of nearest neighbors to determine the class label of unknown record

# *k* NEAREST NEIGHBOR

- Common Distance Metrics:
  - Euclidean distance(continuos distribution)
  $$d(p,q) = \sqrt{\sum(p_i - q_i)^2}$$
  - Hamming distance (overlap metric)
    - **b**at (distance = 1)      **t**o**n**e**d** (distance = 3)
    - **c**at      **r**o**s**e**s**
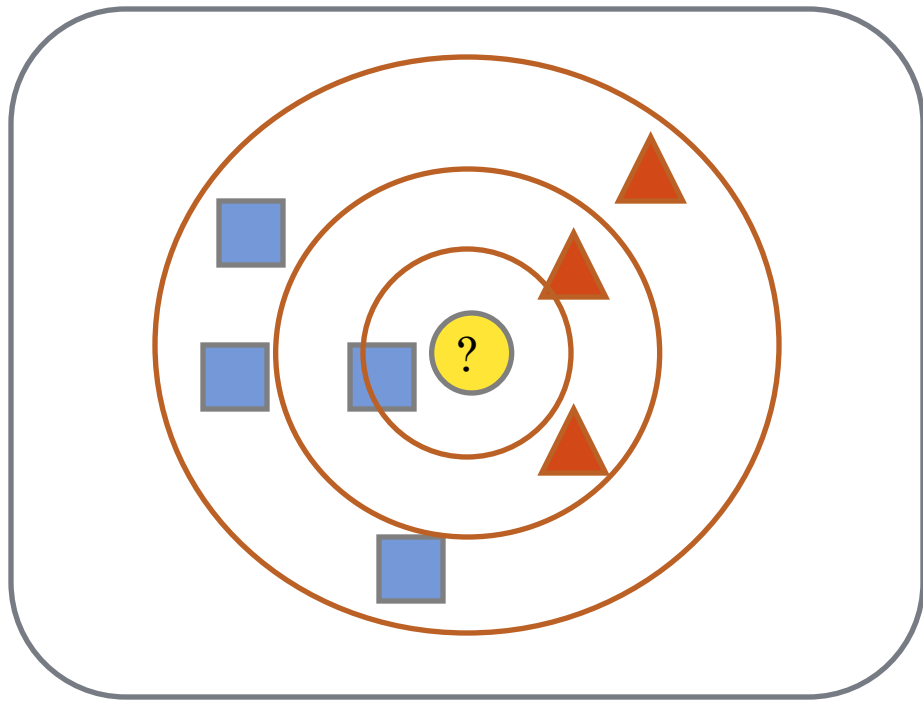  - Discrete Metric(boolean metric)
    - if *x* = *y* then *d(x,y)* = 0. Otherwise, *d(x,y)* = 1

- Determine the class from *k* nearest neighbor list
  - Take the majority vote of class labels among the k-nearest neighbors
  - Weighted factor
    w = 1/d(generalized linear interpolation) or $1/d^2$

# *k* NEAREST NEIGHBOR



- *k* = 1:
  - Belongs to square class

- *k* = 3:
  - Belongs to triangle class

- *k* = 7:
  - Belongs to square class

- Choosing the value of *k*:
  - If *k* is too small, sensitive to noise points
  - If *k* is too large, neighborhood may include points from other classes
  - Choose an odd value for *k*, to eliminate ties

# *k* NEAREST NEIGHBOR

- Accuracy of **all** NN based classification, prediction, or recommendations depends solely on a data model, no matter what specific NN algorithm is used.

- Scaling issues
  - Attributes may have to be scaled to prevent distance measures from being dominated by one of the attributes.
  - Examples
    - Height of a person may vary from 4' to 6'
    - Weight of a person may vary from 100lbs to 300lbs
    - Income of a person may vary from $10k to $500k

- Nearest Neighbor classifiers are lazy learners
  - No pre-constructed models for classification

# ADVANTAGES

- Simple technique that is easily implemented

- Building model is inexpensive

- Extremely flexible classification scheme
  - does not involve preprocessing

- Well suited for
  - Multi-modal classes (classes of multiple forms)
  - Records with multiple class labels

- Asymptotic Error rate at most twice Bayes rate
  - Cover & Hart paper (1967)

- Can sometimes be the best method
  - Michihiro Kuramochi and George Karypis, Gene Classification using Expression Profiles: A Feasibility Study, International Journal on Artificial Intelligence Tools. Vol. 14, No. 4, pp. 641-660, 2005
  - K nearest neighbor outperformed SVM for protein function prediction using expression profiles

# *k* NEAREST NEIGHBOR DISADVANTAGES

- Classifying unknown records are relatively expensive
  - Requires distance computation of k-nearest neighbors
  - Computationally intensive, especially when the size of the training set grows
- Accuracy can be severely degraded by the presence of noisy or irrelevant features

| Height (in cms) | Weight (in kgs) | T Shirt Size |
| --- | --- | --- |
| 158 | 58 | M |
| 158 | 59 | M |
| 158 | 63 | M |
| 160 | 59 | M |
| 160 | 60 | M |
| 163 | 60 | M |
| 163 | 61 | M |
| 160 | 64 | L |
| 163 | 64 | L |
| 165 | 61 | L |
| 165 | 62 | L |
| 165 | 65 | L |
| 168 | 62 | L |
| 168 | 63 | L |
| 168 | 66 | L |
| 170 | 63 | L |
| 170 | 64 | L |
| 170 | 68 | L |

- **Step 1 : Calculate Similarity based on distance function**
- **New customer named 'Monica' has height 161cm and weight 61kg.**

Euclidean :

$$d(x, y) = \sqrt{\sum_{i=1}^{m} (x_i - y_i)^2}$$
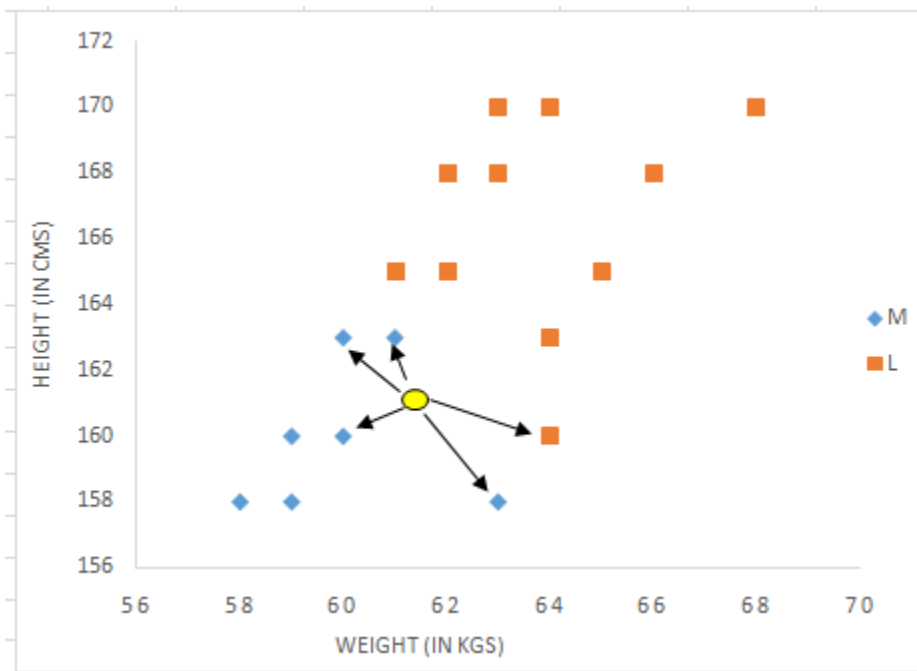
Manhattan / city - block :

$$d(x, y) = \sum_{i=1}^{m} |x_i - y_i|$$

- =SQRT$((161-158)^2+(61-58)^2)$

- Similarly, we will calculate distance of all the training cases with new case and calculates the rank in terms of distance. The smallest distance value will be ranked 1 and considered as nearest neighbor.

# Step 2 : Find K-Nearest Neighbors

fx    =SQRT(($A$21-A6)^2+($B$21-B6)^2)

| | A | B | C | D | E |
|---|---|---|---|---|---|
| 1 | **Height (in cms)** | **Weight (in kgs)** | **T Shirt Size** | **Distance** | |
| 2 | 158 | 58 | M | 4.2 | |
| 3 | 158 | 59 | M | 3.6 | |
| 4 | 158 | 63 | M | 3.6 | |
| 5 | 160 | 59 | **M** | 2.2 | **3** |
| 6 | 160 | 60 | **M** | 1.4 | **1** |
| 7 | 163 | 60 | **M** | 2.2 | **3** |
| 8 | 163 | 61 | **M** | 2.0 | **2** |
| 9 | 160 | 64 | **L** | 3.2 | **5** |
| 10 | 163 | 64 | L | 3.6 | |
| 11 | 165 | 61 | L | 4.0 | |
| 12 | 165 | 62 | L | 4.1 | |
| 13 | 165 | 65 | L | 5.7 | |
| 14 | 168 | 62 | L | 7.1 | |
| 15 | 168 | 63 | L | 7.3 | |
| 16 | 168 | 66 | L | 8.6 | |
| 17 | 170 | 63 | L | 9.2 | |
| 18 | 170 | 64 | L | 9.5 | |
| 19 | 170 | 68 | L | 11.4 | |
| 20 | | | | | |
| 21 | **161** | **61** | | | |

- **Assumptions of KNN**

- **1. Standardization**

    When independent variables in training data are measured in different units, it is important to standardize variables before calculating distance

$$Xs = \frac{X - mean}{s.d.}$$

$$Xs = \frac{X - mean}{max - min}$$

$$Xs = \frac{X - min}{max - min}$$

| | A | B | C | D | E |
|---|---|---|---|---|---|
| 1 | Height (in cms) | Weight (in kgs) | T Shirt Size | Distance | |
| 2 | -1.39 | -1.64 | M | 1.3 | |
| 3 | -1.39 | -1.27 | M | 1.0 | |
| 4 | -1.39 | 0.25 | M | 1.0 | |
| 5 | -0.92 | -1.27 | **M** | 0.8 | **4** |
| 6 | -0.92 | -0.89 | **M** | 0.4 | **1** |
| 7 | -0.23 | -0.89 | **M** | 0.6 | **3** |
| 8 | -0.23 | -0.51 | **M** | 0.5 | **2** |
| 9 | -0.92 | 0.63 | L | 1.2 | |
| 10 | -0.23 | 0.63 | L | 1.2 | |
| 11 | 0.23 | -0.51 | **L** | 0.9 | **5** |
| 12 | 0.23 | -0.13 | L | 1.0 | |
| 13 | 0.23 | 1.01 | L | 1.8 | |
| 14 | 0.92 | -0.13 | L | 1.7 | |
| 15 | 0.92 | 0.25 | L | 1.8 | |
| 16 | 0.92 | 1.39 | L | 2.5 | |
| 17 | 1.39 | 0.25 | L | 2.2 | |
| 18 | 1.39 | 0.63 | L | 2.4 | |
| 19 | 1.39 | 2.15 | L | 3.4 | |
| 20 | | | | | |
| 21 | **-0.7** | **-0.5** | | | |

- **Outlier**
- Low k-value is sensitive to outliers and a higher K-value is more resilient to outliers as it considers more voters to decide prediction

# KNN Exercise

- We have a data from questionnaires survey and objective testing with two attributes(acid durability and strength) to classify whether a special paper tissue is good or not. Here are four training samples .

- Now the factory produces a new paper tissue that pass laboratory test

| X1= Acid durability (seconds) | X2=Strength (kg/m2) | Y=Classification |
|---|---|---|
| 7 | 7 | Bad |
| 7 | 4 | Bad |
| 3 | 4 | Good |
| 1 | 4 | Good |

with X1=3 and X2=7, without another expensive survey can we guess what the classification of new tissue is

Step 1:Determine parameter k= number of nearest neighbors.

Suppose K=3

Step 2: Calculate the distance between query instance and all the training samples.

| X1= Acid durability (seconds) | X2=Strength (kg/m2) | Distance |
|---|---|---|
| 7 | 7 | 4 |
| 7 | 4 | 5 |
| 3 | 4 | 3 |
| 1 | 4 | 3.6 |

- Step3: Sort the distance and determine the nearest neighbor based on the Kth minimum distance.

| X1= Acid durability (seconds) | X2=Strong th (kg/m2) | Distance | Rank | is it included in 3NN? |
|---|---|---|---|---|
| 7 | 7 | 4 | 3 | YES |
| 7 | 4 | 5 | 4 | No |
| 3 | 4 | 3 | 1 | YES |
| 1 | 4 | 3.6 | 2 | YES |

- Step 4: Use simple majority of the category of nearest neighbour as the prediction value for query instance.

- Therefore we have 2 good and one bad . Since 2>1.

 a new paper tissue that pass laboratory test with $X_1=3$ and $X_2=7$ is included in Good category