

## Unit 2: Dimensionality Reduction

By Prof. Ranjitha U.N

### Introduction

In machine learning as the dimensionality of the data rises, the amount of data required to provide a reliable analysis grows exponentially. Bellman referred to this phenomenon as the "curse of dimensionality".

### Cases:

1. A motorbike rider in racing competitions. His position and movement gets measured by GPS sensor on bike, gyro meters, multiple video feeds and his smart watch. Because of respective errors in recording, the data would not be exactly same. However, there is very little incremental information on position gained from putting these additional sources.

Now assume that an analyst sits with all this data to analyze the racing strategy of the biker – he/ she would have a lot of variables / dimensions which are similar and of little (or no) incremental value. This is the problem of high unwanted dimensions and needs a treatment of dimension reduction

2. **Image processing.** Facebook application – "Which Celebrity Do You Look Like?".

To identify the matched celebrity image, we use pixel data and each pixel is equivalent to one dimension. In every image, there are high number of pixels i.e. high number of dimensions. And every dimension is important here. You can't omit dimensions randomly to make better sense of your overall data set. In such cases, dimension reduction techniques help you to find the significant dimension(s) using various method(s).

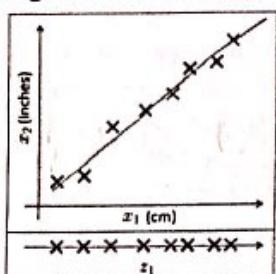


Figure 1: 2D converted to one dimension

3. There are many sources of data that can be viewed as a large matrix. Web can be represented as a transition matrix.

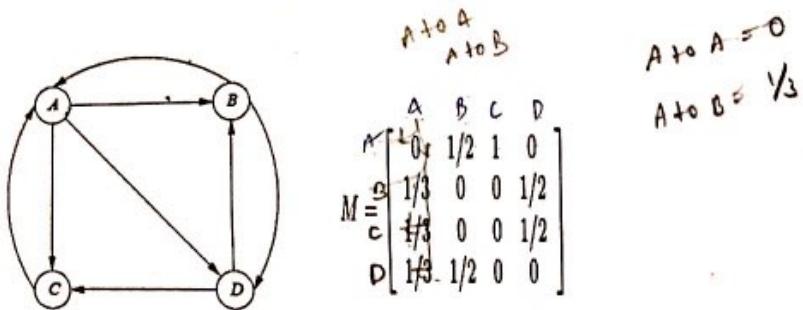


Fig: Hypothetical example of web

 While Google was not the first search engine, it was the first, able to defeat the spammers, PageRank was established as an essential technique for a search engine to defeat link spam

Social network can also be represented as a transition matrix.

In many of these matrix applications, the matrix can be summarized by finding “narrower” matrices (linear algebra, dot product and eigen values) that in some sense are close to the original. These narrow matrices have only a small number of rows or a small number of columns, and therefore can be used much more efficiently than can the original large matrix. The process of finding these narrow matrices is called dimensionality reduction.

**Definition :** Dimension Reduction refers to the process of converting a set of data having vast dimensions into data with lesser dimensions ensuring that it conveys similar information concisely in turn to minimize the complexity of any classifier or regressor

**Reason /objective:**

1. In most learning algorithms, the complexity depends on the number of input dimensions,  $d$ , as well as on the size of the data sample,  $N$ , and for reduced memory and computation, we are interested in reducing the dimensionality of the problem. Decreasing  $d$  also decreases the complexity of the inference algorithm during testing.
2. When an input is decided to be unnecessary, we save the cost of extracting it.
3. Simpler models are more robust on small datasets (exhibit less variance)
4. When data can be explained with fewer features, we get a better idea about the process that underlies the data and this allows knowledge extraction.
5. When data can be represented in a few dimensions without loss of information, it can be plotted and analyzed visually for structure and outliers

**Methods:** There are two main methods for reducing dimensionality: feature selection and feature extraction.(can be supervised or unsupervised)

In feature selection, finding  $k$  of the  $d$  dimensions that give us the most information and we discard the other  $(d - k)$  dimensions. Ex. Random forest , decision tree

In feature extraction, we are interested in finding a new set of  $k$  dimensions that are combinations of the original  $d$  dimensions ex: Principal Components Analysis (PCA) and Linear Discriminant Analysis (LDA),

unsupervised  
PCA & LDA  
Projection  
supervised  
(FA, MDS)

### Subset selection (feature subset selection):

In subset selection, we are interested in finding the best subset of the set of features. The best subset contains the least number of dimensions that most contribute to accuracy .

For  $d$  dimesions ,there are  $2^d$  possible subsets of  $d$  variables, but we cannot test for all of them.

There are two approaches in selecting the subset:

Forward selection

Backward selection

Forward selection: - we start with no variable and add them one by one, at each step adding the one that decreases the error the most, until any further addition does not decrease the error

Sequential forward selection (algorithm)

$E(F) \rightarrow$  error on  $\{F\}$   
increasing classification error

→ Start with no features:  $F = \emptyset$ .  $F$  , is a feature set of dimension  $x_i : i = 1 \dots d$

→ At each step, for all possible  $x_i$  , train our model on the training set and calculate  $E(F \cup x_i)$  on the validation set.

→ choose that input  $x_j$  that causes the least error

→  $j = \arg \min_i E(F \cup x_i)$  and add  $x_j$  to  $F$  if  $E(F \cup x_j) < E(F)$

Drawback: because this algorithm is greedy and adds attributes one by one, it may not be able to detect the effectiveness of  $x_i$  and  $x_j$  , as both together can reduce error

Backward selection: - we start with all variables and remove them one by one, at each step removing the one that decreases the error the most, until any further removal increases the error significantly.

### Dimensionality reduction

Prof . Ranjitha U.N, School of C&IT, RU

Forward selection this dataset has 4 classes & 3 features. Accuracy 0.76, 0.52, 0.92 & 0.94  
1. classifier Nearest mean → 1st 4 features accuracy 0.76, 0.52, 0.92 & 0.94  
in 3d also 0.94  
 $d + (d-1) + (d-2) \dots + (d-k)$   
 $O(d^2)$

Let  $F$  ---- feature set of input dimensions,  $x_i, i = 1 \dots d$   
 $E(F)$  ----- denotes the error incurred on the validation sample (distinct from training sample)  
the error is either the mean square error or misclassification error

#### Sequential backward selection

\*Start with  $F$  containing all features

\*remove one attribute from  $F$  as the one that causes the least error i.e.

$$j^* = \arg \min_i E(F - x_i)$$

and we remove  $x_j$  from  $F$  if  $E(F - x_j) < E(F)$

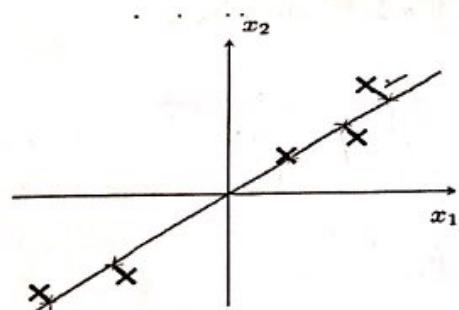
Drawback: training a system with more features is more costly than training a system with fewer features .

Note: forward search may be preferable especially if we expect many useless features.  
In supervised learning, output is considered to calculate the error during subset selection.

#### **Feature Extraction:**

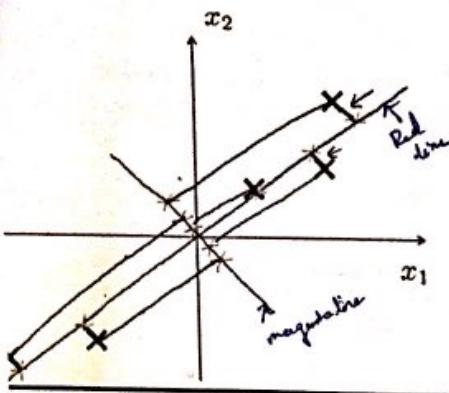
In applications like face recognition, removal of feature (pixel), may lead to loss of information as group of pixels together carry information. Such cases deal with feature extraction than feature removal.

PCA is an unsupervised feature extraction algorithm



Reduce 2 dimension to one dimension (similarly n dimension to k dimension).  
To do so find a direction (vector) on to which the data is projected so as to minimize the projection error (blue lines in fig)  
Note: PCA is not linear regression,

Priffid  $\Rightarrow$  more variance  
less error



The criterion in PCA is to maximize variance .(the line in red is better when compared to magenta line which has more projection error well as less variance).

It is always better to maintain good distance between the points i.e. difference between the sample points becomes most apparent

In general the projection of  $x$  on the direction of  $w$  is  
$$z = w^T x$$

## Factor analysis

Preread : There has been significant controversy in the field over differences between the two techniques ,exploratory factor analysis and principal components analysis ). PCA can be considered as a more basic version of exploratory factor analysis (EFA) that was developed in the early days prior to the advent of high-speed computers.

Both PCA and factor analysis aim to reduce the dimensionality of a set of data.

**Factor analysis :-** is clearly designed with the objective to identify certain unobservable factors from the observed variables, whereas PCA does not directly address this objective

There are two uses of factor analysis:

- ① It can be used for knowledge extraction ,  $\rightarrow$  grouping of highly correlated variables under one factor and others in other factor through subset creation
- ② It can also be used for dimensionality reduction when  $k < d$

FOR DIMENSIONALITY REDUCTION

In factor analysis (FA), we assume that there is a set of unobservable, latent factors  $z_j$ ,  $j = 1, \dots, k$ , which when acting in combination generate  $x$ .

When we are interested in dimensionality reduction, we need to be able to find the factor scores,  $z_j$ , from  $x_i$ . We want to find the loadings  $w_{ji}$

*Explanation of factor loading*

Dimensionality reduction

Prof . Ranjitha U.N, School of C&IT, RU

Note: Large no. of variables into fewer no. of factors.

articulators, namely  
to shape the air as  
TJ

\* Latent factors  $z_i$ ,  $i=1 \dots K$  -> Latent variables act  
combination generate  $\underline{x}$  (i/p)

Variables	Factor 1
Income	0.65
Education	0.59
Occupation	0.48
House value	0.38
Number of public parks in neighborhood	0.13
Number of violent crimes per year in neighborhood	0.23

Note: The relationship of each variable to the underlying factor is expressed by the so-called factor loading.

The correlation between the variable and the factor is given by the score. Hence the FA deals with the correlation where PCA deals with variance. Here, latent variable is socio-economic status

We want to find the loadings  $w_{ji}$ , such that

$$z_j = \sum_{i=1}^d w_{ji} x_i + \epsilon_i, j = 1, \dots, k \quad (\text{Eq. 1}), \quad \text{where } x_i \text{ are centered to have mean 0}$$

In vector form, for observation T, it can be written as

$$z^t = W^T x^t + \epsilon, \forall t = 1, \dots, N$$

This is a linear model with  $d$  inputs and  $k$  outputs. Its transpose can be written as

$$(z^t)^T = (x^t)^T W + \epsilon^T, \forall t = 1, \dots, N$$

For a sample of  $N$  observations, we have

$$Z = XW + \Sigma \quad \Sigma \text{ is } N \times k \text{ of zero-mean noise.}$$

— (2)

$W$  can be found as

$$\begin{aligned} W &= (N-1)(X^T X)^{-1} \frac{X^T Z}{N-1} \\ &= \left( \frac{X^T X}{N-1} \right)^{-1} \frac{X^T Z}{N-1} \quad \text{Substituting the values for } W \\ &= S^{-1} V \end{aligned}$$

— (3)

$$Z = XW = XS^{-1}V$$

### Comparison between PCA and FA

FA offers no advantage over PCA except the interpretability of factors allowing the identification of common causes, a simple explanation, and knowledge extraction. For example, in the context of speech recognition,  $x$  corresponds to the acoustic signal,

Dimensionality reduction

Prof. Ranjitha U.N, School of C&IT, RU

articulators, namely, jaw, tongue, velum, lips, and mouth, which are positioned appropriately to shape the air as it comes out of the lungs and generate the speech sound.

If speech signal could be transformed to these attributes, then recognising would be easy.

### Linear Discriminant Analysis

Linear Discriminant Analysis is a supervised method for dimensionality reduction for classification problems. Now considering two classes, then generalizing to  $K > 2$  classes.

- \* Considering samples from two classes  $C_1$  and  $C_2$ , a new direction defined by a vector  $w$  is found such that when data are projected onto  $w$ , the examples from the two classes are well separated.

$$z = w^T x \Rightarrow \text{projection of } x \text{ onto } w. [\Rightarrow \text{Reduction from } d \text{ to 1}]$$

- \*  $m_1 \rightarrow$  means of sample from  $C_1$  before projection

$$m'_1 \rightarrow \text{''} \quad \text{''} \quad \text{''} \quad C_1 \text{ after projection}$$

$$m_1 \in \mathbb{R}^d \quad \text{and} \quad m'_1 \in \mathbb{R}$$

- \* For a given sample  $x = \{x^t, \gamma^t\}$

$$\gamma^t = 1 \quad \text{if } x^t \in C_1 ; \quad \gamma^t = 0 \quad \text{if } x^t \in C_2$$

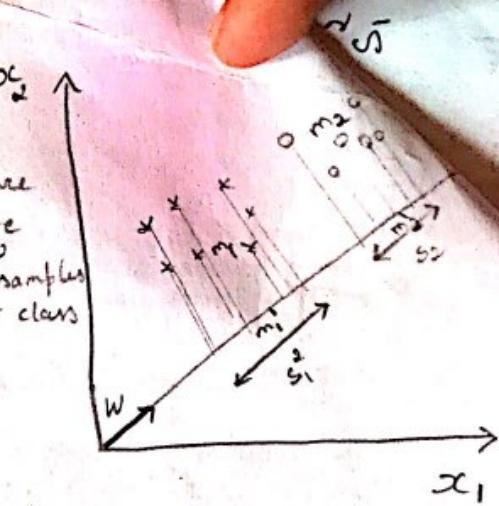
$$m'_1 = \frac{\sum_t w^T x^t \gamma^t}{\sum_t \gamma^t} = w^T m_1$$

$$m'_2 = \frac{\sum_t w^T x^t (1 - \gamma^t)}{\sum_t (1 - \gamma^t)} = w^T m_2$$

The scatter of samples from  $C_1$  and  $C_2$  after projection are

$$S_1^2 = \sum_t (\omega^T x^t - m_1^t)^2 \quad \begin{matrix} \text{Scatter = Variance} \\ \text{sum of square} \\ \text{differences b/w} \\ \text{projected samples} \\ \text{and their class} \\ \text{mean.} \end{matrix}$$

$$S_2^2 = \sum_t (\omega^T x^t - m_2^t)^2$$



For the two classes to be well separated, the means need to be as far as possible and the examples of classes be scattered in as small region as possible.

i.e (from fig) we need  $|m_1^t - m_2^t|$  to be large and  $S_1^2 + S_2^2$  to be small  $\Rightarrow$  (Fisher's linear discriminant)  $w$  that maximizes below eq

$$\mathcal{J}(w) = \frac{(m_1^t - m_2^t)^2}{S_1^2 + S_2^2} \quad \text{--- (1)}$$

Rewriting the numerator, we get

$$\begin{aligned} (m_1^t - m_2^t)^2 &= (\omega^T m_1 - \omega^T m_2)^2 \\ &= \underbrace{\omega^T (m_1 - m_2)}_{S_B} \underbrace{(m_1 - m_2)^T}_{(m_1 - m_2)^T} \omega \end{aligned}$$

Note:-  
scatter matrix  
is a statistic used  
to make a estimate  
of covariance matrix

where,  $S_B = (m_1 - m_2)(m_1 - m_2)^T$  is the between-class scatter matrix.

The denominator of eq (1) (which the sum of scatter) can be written as

$$S_1^2 = \sum_t (\omega^T x^t - m_1)^2 \gamma_t^t$$

$$= \sum_t \{ \omega^T (x^t - m_1) (x^t - m_1)^T \omega \} \gamma_t^t$$

$$= \omega^T S_1 \omega$$

where  $S_1 = \sum_t (x^t - m_1) \gamma_t^t (x^t - m_1)^T$  is the within class scatter matrix

Applying the same (above) analysis to  $S_2$

we have  $S_1^2 + S_2^2 = \omega^T S_w \omega$ .  $\frac{S_w}{\downarrow} = S_1^2 + S_2^2$

Total within class scatter matrix.

Note:  $\frac{S_1^2 + S_2^2}{\text{Total no. of samples}} = \text{Variance of pooled data}$

When  $K > 2$  classes

matrix  $\omega$  needs to be found such that

$$z = \omega^T \otimes x. \quad z \rightarrow K \text{ dimensional}$$

$$\omega \rightarrow d \times K$$

within-class matrix for  $c_i$  is

$$S_i = \sum_t \gamma_i^t (x^t - m_i) (x^t - m_i)^T$$

Total within-class scatter is

$$S_w = \sum_{i=1}^K S_i$$

The between-class scatter matrix is

$$S_B = \sum_{i=1}^K N_i (m_i - \bar{m}) (m_i - \bar{m})^T \quad N_i = \sum_t \gamma_i^t$$

For the analysis, we need

1. the between-class scatter matrix  $w^T S_B w$  to be large
2. the within-class scatter matrix  $w^T S_w w$  to be small

Note: For a scatter matrix,

Measure of the spread = determinant of the matrix

Determinant of a matrix = product of eigenvalues  
and

Eigenvalues gives variance along its eigenvector

∴ The interest is in the matrix  $W$  that maximizes

$$\xi(w) = \left| \frac{w^T S_B w}{w^T S_w w} \right|$$

The largest eigenvectors of  $S_w^{-1} S_B$  are the solution.

## Multidimensional scaling

- It is also a dimensionality reduction algorithm like PCA and FA, but the approach taken for dim reduction varies.
- It focusses on the distance between the related item (data point) and not on their spatial coordinates (positions).
- Implementation of MDS, can be
  - i) Metric ie Classical Multidimensional Scaling (CMDS) because it tries to reproduce the original metric or distances when the data is in metric form
  - ii) Non-Metric Multidimensional Scaling (NMMDS), this method produces a map which tries to reproduce ranks when the data is in the form of ranking or ordinal data

Consider the following problem: looking at a map showing a number of cities, one is interested in the distances between them. These distances are easily obtained by measuring them using a ruler. Apart from that, a mathematical solution is available: knowing the coordinates  $x$  and  $y$ , the Euclidean distance between two cities  $a$  and  $b$  is defined by

$$d_{ab} = \sqrt{(x_a - x_b)^2 + (y_a - y_b)^2}$$

Now consider the inverse problem: having only the distances, is it possible to obtain the map? Classical MDS, addresses this problem.

## Steps of a classical MDS algorithm

**NOTE:** Classical MDS algorithms typically involve some linear algebra.

The classical MDS algorithm rests on the fact that the coordinate matrix  $X$  can be derived by eigenvalue decomposition from the scalar product matrix  $B = XX'$ .

The problem of constructing  $B$  from the proximity matrix  $P$  is solved by the method called double centering Where the squared proximity matrix is multiplied with the matrix  $J$

$$J = I - \frac{1}{n} \mathbf{1} \mathbf{1}' , \text{ where } n \text{ is the number of objects.}$$

The following steps summarize the algorithm of classical MDS

- Set up the matrix of squared  $P^{(2)} = [p^2]$  proximities
- Apply double  $B = -\frac{1}{2}JP^{(2)}J$ , centering
- Extract the  $m$  largest positive eigenvalues  $\lambda_1, \dots, \lambda_m$  of  $B$  and the corresponding  $m$  Eigen vectors  $e_1, \dots, e_m$
- A  $m$ -dimensional spatial configuration of the  $n$  objects is derived from the coordinate matrix  $X = E_m \Lambda_m^{1/2}$

Where  $E_m$  is the matrix of  $m$  eigenvectors and  $\Lambda_m$  is the diagonal matrix of  $m$  eigenvalues of  $B$ , respectively

Illustration of MDS:

Assume that we have a proximity matrix that has been obtained from the distance measured between city  $c_1$  &  $c_2$  on a map cities  $c_1, c_2, c_3$  and  $c_4$  on a map

$$P = \begin{bmatrix} c_1 \\ c_2 \\ c_3 \\ c_4 \end{bmatrix}$$

and  $c_4$  on a map

{in(mm)}

The proximity matrix of the cities is given by

$$P = \begin{bmatrix} c_1 & 0 & 93 & 82 & 133 \\ c_2 & 93 & 0 & 52 & 60 \\ c_3 & 82 & 52 & 0 & 111 \\ c_4 & 133 & 60 & 111 & 0 \end{bmatrix}$$

Step 1:  $P^2 = \begin{bmatrix} 0 & 8649 & 6724 & 17689 \\ 8649 & 0 & 2704 & 3600 \\ 6724 & 2704 & 0 & 12321 \\ 17689 & 3600 & 12321 & 0 \end{bmatrix}$

Since  $n=4$  objects, the matrix  $J$  is calculated by

$$J = \begin{bmatrix} 1 & 0 & 0 & 0 \\ 0 & 1 & 0 & 0 \\ 0 & 0 & 1 & 0 \\ 0 & 0 & 0 & 1 \end{bmatrix} - 0.25 \times \begin{bmatrix} 1 & 1 & 1 & 1 \\ 1 & 1 & 1 & 1 \\ 1 & 1 & 1 & 1 \\ 1 & 1 & 1 & 1 \end{bmatrix} = \begin{bmatrix} 0.75 & -0.25 & -0.25 & -0.25 \\ -0.25 & 0.75 & -0.25 & -0.25 \\ -0.25 & -0.25 & 0.75 & -0.25 \\ -0.25 & -0.25 & -0.25 & 0.75 \end{bmatrix}$$

unit matrix

diagonal unit

Obtaining double centered matrix  $B$

$$B = -\frac{1}{2} J P^2 J = \begin{bmatrix} 5035.06 & -1553.06 & 258.93 & -3740.9 \\ -1553.06 & 507.8 & 5.312 & 1039 \\ 258.9 & 5.31 & 2206 & -247.06 \\ -3740.9 & 1039.9 & -2471.6 & 5172.06 \end{bmatrix}$$

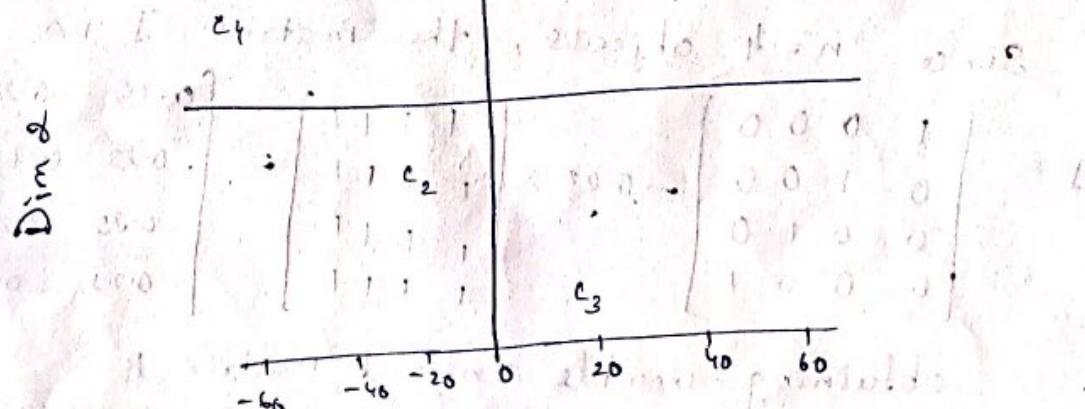
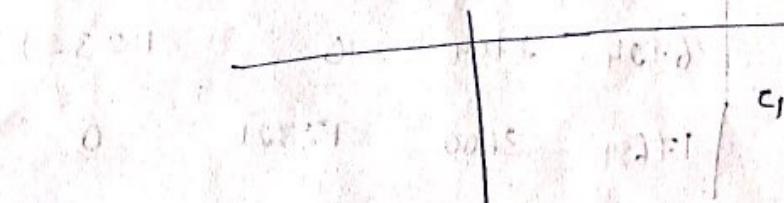
For a two-dimensional representation of the four cities, the first two largest eigenvalues and the corresponding eigenvectors of  $B$  have to be extracted.

$$\lambda_1 = 9724.16 \quad \lambda_2 = 3160.98$$

$$e_1 = \begin{bmatrix} -0.63 \\ 0.18 \\ -0.25 \\ 0.704 \end{bmatrix}, \quad e_2 = \begin{bmatrix} -0.586 \\ 0.214 \\ 0.706 \\ -0.334 \end{bmatrix}$$

$m$ -dimensional spatial configuration gives the co-ordinates of the cities are obtained by multiplying eigenvalues and vectors.  $\therefore X = E_m \lambda_m^{1/2}$

$$X = \begin{bmatrix} -0.63 & -0.58 \\ 0.18 & 0.214 \\ -0.25 & 0.706 \\ 0.704 & -0.334 \end{bmatrix} \begin{bmatrix} \sqrt{9724.16} & 0 \\ 0 & \sqrt{3160.98} \end{bmatrix} = \begin{bmatrix} c_1 & -32.9 \\ c_2 & 12.02 \\ c_3 & 39.71 \\ c_4 & -18.76 \end{bmatrix}$$



## K-Nearest-Neighbors classifiers

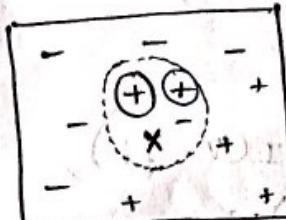
Prof. Ranjitha U-N

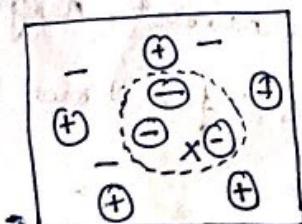
- \* Nearest-neighbor classifiers are based on comparing a given test tuple with training tuples that are similar to it.
- \* All of the training tuples are stored in an n-dimensional pattern space. Each tuple represents a point in n-dimensional space.
- \* When an unknown tuple is given, a K-nearest-neighbor classifier searches the pattern space for the K training tuples that are closest to the unknown tuple. Thus the K nearest neighbors of the unknown tuple is found.

\* Closeness is defined in terms of a distance metric such as Euclidean distance.

$$\text{For 2 tuples } x_1 = (x_{11}, x_{12}, \dots, x_{1n}) \text{ & } x_2 = (x_{21}, x_{22}, \dots, x_{2n})$$

$$\text{dist}(x_1, x_2) = \sqrt{\sum_{i=1}^n (x_{1i} - x_{2i})^2}$$

- \*   $\rightarrow$  2 nearest neighbor  
Label of the new sample 'x' will be the label of '+' samples as they are the majority vote.



$$K = 3$$

## Algorithm:

1. Let  $K$  be the number of nearest neighbors  
and  $D$  be the set of training examples.
2. for each test example  $z = (x', y')$  do
3. compute  $d(x', x)$ , the distance between  $z$  and  
every example  $(x, y) \in D$ .
4. Select  $D_2 \subseteq D$ , the set of  $K$  closest training  
examples to  $z$ .
5.  $y' = \arg \max_v \sum_{(x_i, y_i) \in D_2} I(v = y_i)$
6. end for.

## Summarizing the algorithm

Algorithm computes the distance between  
each test example  $z = (x', y')$  and all the  
training example  $(x, y) \in D$  to determine its  
nearest neighbor list  $D_2$ .

Once the nearest neighbor list is obtained,  
the test example is classified based on majority  
class of its nearest neighbors.

$$\text{Majority Voting: } y' = \arg \max_v \sum_{(x_i, y_i) \in D_2} I(v = y_i),$$

$v \rightarrow$  class label

$y_i \rightarrow$  class label for one of the nearest neighbors.

$I(\cdot) \rightarrow$  Indicator function returns 1 for true & 0  
otherwise.

Apply KNN on the dataset of the company which produces tissues for laboratory. Predict the acceptability of the new type of tissue with acid durability being 3 & strength being 7. for  $K=2$ .

Survey is costly

Solu.

Name	Acid Durability	Strength	Class
Type-1	7	7	Bad
Type-2	7	4	Bad
Type-3	3	4	Good
Type-4	1	4	Good

Given data  
new

Distance

$$D(\text{Type 5, Type 1}) = \sqrt{(7-3)^2 + (7-7)^2} = 4 \Rightarrow \text{Bad}$$

$$D(\text{", Type 2}) = \sqrt{(7-3)^2 + (4-7)^2} = 4 \Rightarrow \text{Bad}$$

$$D(\text{Type 5, Type 3}) = \sqrt{(7-3)^2 + (4-7)^2} = 4 \Rightarrow \text{Good}$$

$$D(\text{Type 5, Type 4}) = \sqrt{(7-1)^2 + (4-7)^2} = 2 \Rightarrow \text{Good}$$

For  $K=2$  in considering 2 immediate neighbour

belongs to class "Good".

Type-5 belongs to class Good.

Prof. Rayitha U.N

\* Apply KNN algorithm and predict the type of food to which Tomato belongs to consider the attribute (sweet = 6 & crunch = 4).

<u>Solu</u>	Ingredient	sweet	crunch	food type
	apple	8	5	fruit
	green bean	3	7	vegetable
	Nuts	3	6	protein
	Orange	7	3	fruit

$$D(\text{tomato}, \text{apple}) = \sqrt{(8-6)^2 + (4-5)^2} = 2.2$$

$$D(\text{,,}, \text{green bean}) = \sqrt{(6-3)^2 + (4-7)^2} = 4.2$$

$$D(\text{,,}, \text{Nuts}) = \sqrt{(6-3)^2 + (4-6)^2} = \sqrt{13} = 3.6$$

$$D(\text{,,}, \text{orange}) = \sqrt{(6-7)^2 + (4-3)^2} = 1.4$$

The min distance is 1.4. Hence Tomato will belong to food type of orange ie fruit.

## Pearson Correlation Coefficient

- \* It is also known as the "product moment correlation coefficient" PMCC or simply correlation.
- \* Pearson's correlations are suitable only for metric variables.
- \* Correlation is a technique for investigating the relationship between two quantitative, continuous variables. (ex: age v/s blood pressure).
- \* Pearson's coefficient ( $\rho$ ) is a measure of the strength of the association between the two variables
- \* The value of the coeff is varies from -1 to 1
- \* If  $\rho = -1$ , data lie on a perfect straight line with a negative slope
- \* If  $\rho = 0$ , no linear relation b/w the variables

If  $\gamma = +1$ , data lie on  
a perfect straight line  
with a positive slope

\* positive correlation ( $\gamma = +1$ ) indicates that both variables increases or decreases together.

Negative correlation ( $\gamma = -1$ ) indicates that as one variable increases, the other decreases & vice versa

$$\text{Pearson's coefficient } \gamma = \frac{1}{n-1} \left( \sum (x - \bar{x})(y - \bar{y}) \right)$$

$$\sqrt{\frac{1}{n-1} \sum (x - \bar{x})^2} \cdot \sqrt{\frac{1}{n-1} \sum (y - \bar{y})^2}$$

Simplified as

$$\gamma = \frac{\sum xy - \frac{1}{n} \sum x \sum y}{\sqrt{\frac{\sum x^2 - (\sum x)^2}{n}} \sqrt{\frac{\sum y^2 - (\sum y)^2}{n}}}$$

Compute Pearson's coefficient of correlation between advertisement cost and sales as per the data given below.

Advertisement cost in 1000's	39	65	62	90	82	75	25	98	36	78
	47	53	58	86	62	68	60	91	51	84
Sales in lakhs										X

Solve From the data,  $n = 10$

$$\sum x = 650 \quad \sum y = 660 \quad \sum xy = 45604$$

$$\sum x^2 = 47648 \quad \sum y^2 = 45784$$

$$r = \frac{45604 - \frac{(650)(660)}{10}}{\sqrt{47648 - \frac{(650)^2}{10}} \sqrt{45784 - \frac{(660)^2}{10}}}$$

$$r = \frac{45604 - 42900}{(73.47)(47.1)} = 0.7804$$

Correlation coeff is positively correlated.

Decision tree Learning :-

- \* It is a practical method for inductive inference.
- \* This method search a completely expressive hypothesis space and avoid the difficulties of restricted hypothesis spaces.
- \* Widely used D.T algorithms are ID 3, ASSISTANT etc
- \* It represents a flow-chart kind of model which can also be represented as sets of if-then rules.
- \* Applications of Decision tree algorithm varies from diagnose of medical case to learning to assess credit risk of loan applicants.

Decision tree Representation :-

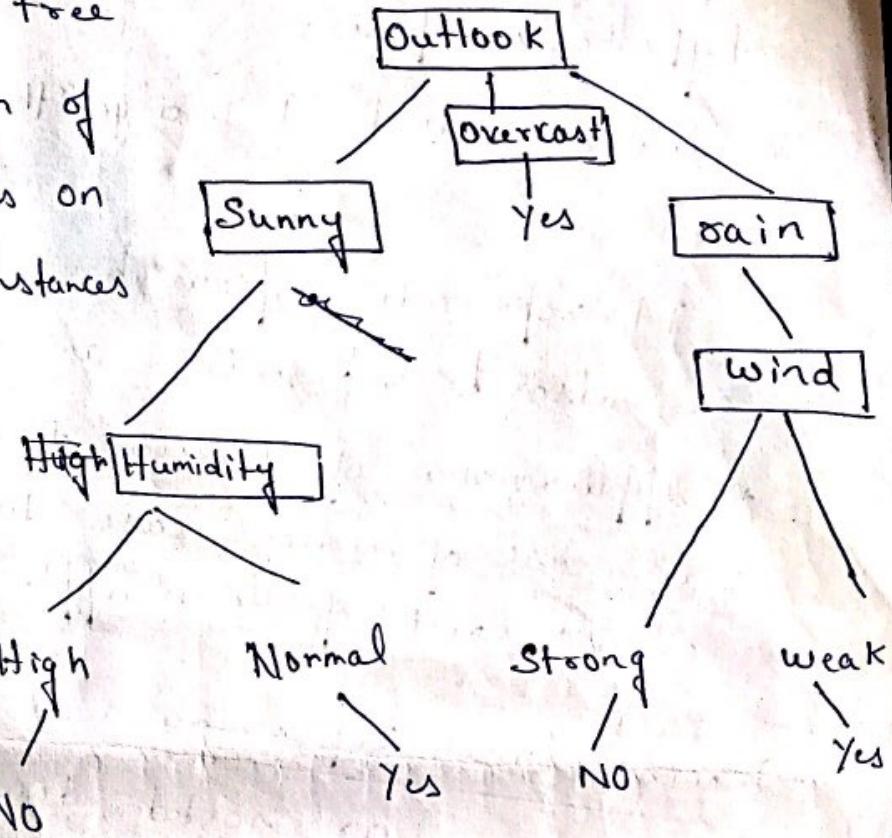
- Instances are sorted down the tree from the root to some leaf node to achieve classification.
- Each node in the tree specifies a test of some attribute of the instance and each branch from that node corresponds to one of the possible values for this attribute.

## Illustration of decision tree :-

Classification of Saturday mornings according to whether they are suitable for playing tennis.

In general, decision tree represents a disjunction of conjunctions of constraints on the attribute values of instances

The decision tree shown corresponds to the expression



$$(\text{outlook} = \text{sunny} \wedge \text{Humidity} = \text{Normal})$$

$$\vee (\text{outlook} = \text{Overcast})$$

$$\vee (\text{outlook} = \text{Rain} \wedge \text{wind} = \text{weak})$$

## Appropriate problems for Decision Tree Learning

Decision tree learning is generally best suited to problems with the following characteristics:

1. Instances are represented by attribute-value pairs.  
Easiest situation for decision tree learning is when each attribute takes on a few possible values.  
Ex: If temperature is a variable (attribute).  
possible values would be hot, mild, cold.
2. The target function has discrete output values.  
Ex: Obj with Boolean values  
More than two values  
some real valued values.
3. Requirement of disjunctive descriptions
4. Decision tree learning methods are robust to errors being in training examples and also the errors in attribute values.
5. The training data may contain missing attribute values

$$G.P. = -(P_0 \log_2 P_0 + P_1 \log_2 P_1)$$

and we want to find the G.P. for class 1

Which attribute is the best classifier?

For classification using decision tree, we need to select an attribute that is most which gives the best result.

In order to achieve best result, a statistical property called Information gain is considered.

ID3 uses this information gain measure to select among the candidate attributes at each step while growing the tree.

\* Information gain is precisely calculated using entropy

$$\text{Entropy}(S) = -P_+ \log_2 P_+ - P_- \log_2 P_-$$

S relates to boolean classification

$P_+$  → Portion of positive examples in S.

$P_-$  → Portion of negative examples in S

\* if  $P_+ = 9/14$  &  $P_- = 5/14$   
then

$$\text{Entropy}([9+, 5-]) = -(9/14) \log_2 (9/14) + 5/14 \log_2 (5/14)$$
$$= 0.413 + 0.159 + 0.530$$
$$= 0.940$$

Note: If all members of S belong to the same class  
ie  $P_+ = 1$ , then  $P_- = 0$  then  $S = 0$

$$\text{i.e. } S = -1 \cdot \log_2 1/14 - 0 \cdot \log_2 0$$
$$= 0$$

$S = +$

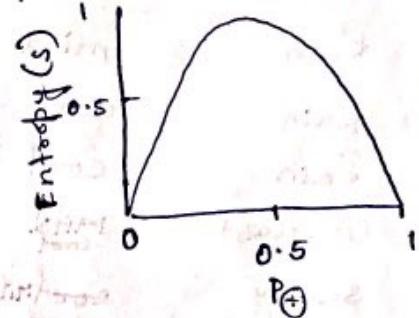
Note: If the sample contains equal no. of  $P_+$  &  $P_-$  then entropy  $S = 1$

If collection contains unequal no. of positive & negative examples, the entropy is between 0 & 1  
Represented as

Entropy (from Information theory)

(perspective) is the min. no. of bits of info needed to encode the classification. i.e.

if  $P_+$  is 1; Receiver knows drawn sample is true  $\Rightarrow$  no msg need to be sent ie  $S = 0$



Note: When the target attribute can take on  $c$  different values, then the entropy of  $s$  is given by

$$\text{Entropy}(s) = \sum_{i=1}^c -P_i \log_2 P_i$$

The measure, information gain is the expected reduction in entropy caused by partitioning the examples according to the attribute. i.e.

$$\text{Gain}(s, A) = \text{Entropy}(s) - \underbrace{\sum_{v \in \text{values}(A)} \frac{|S_v|}{|s|} \text{Entropy}(S_v)}_{\text{expected entropy}}$$

$A \rightarrow$  attribute

$s \rightarrow$  collection of samples

$\text{values}(A) \rightarrow$  all possible values of attribute  $A$ .

Day	Outlook	Temperature	Humidity	Wind	Play tennis
1	Sunny	Hot	High	Weak	No
2	Sunny	Hot	High	Strong	No
3	Overcast	Hot	High	Weak	Yes
4	Rain	Mild	High	Weak	Yes
5	Rain	Cool	Normal	Weak	Yes
6	Rain	Cool	Normal	Strong	No
7	Overcast	Mild cool	Normal	Strong	Yes
8	Sunny	Cool Mild	High	Weak	No
9	Sunny	Cool	Normal	Weak	Yes
10	Rain	Mild	Normal	Weak	Yes
11	Sunny	Mild	Normal	Strong	Yes
12	Overcast	Mild	High	Strong	Yes
13	Overcast	Hot	Normal	Weak	Yes
14	Rain	Mild	High	Strong	No

Calculating gain for each attribute (outlook, temp, humidity)

Information gain for outlook (sunny, overcast, rain)

For sunny :-  $P = 2$   
 $N = 3$

$$\text{gain}(S, \text{sunny}) = -\frac{2}{5} \log_2 \frac{2}{5} - \frac{3}{5} \log_2 \frac{3}{5}$$
$$= 0.529 + 0.442$$
$$= 0.971$$

for outlook = overcast

$$P = 4 \quad \& \quad N = 0$$

$$\text{Entropy}_{\text{overcast}} = \frac{4}{4} \log_2 \frac{4}{4} = 0$$

for outlook = rain

$$P = 3 \quad \& \quad N = 2$$

$$\text{Entropy}_{\text{rain}} = 0.971$$

Entropy of outlook

$$E(\text{outlook}) = \frac{5}{14} I(2, 5) + \frac{4}{14} I(4, 0) + \frac{5}{14} I(3, 2)$$

$$= \frac{5}{14} \times 0.971 + \frac{5}{14} \times 0.971 = 0.694$$

$$\text{gain}(\text{outlook}) = I(P, n) - E(\text{outlook})$$

$$= 0.940 - 0.694$$

$$= 0.246$$

Information gain for Humidity (High, Normal)

For high:  $P = 3$   
 $N = 4$

$$\begin{aligned}E_{\text{high}} &= -\frac{3}{7} \log_2 \left(\frac{3}{7}\right) - \frac{4}{7} \log_2 \left(\frac{4}{7}\right) \\&= -\frac{3}{7} \log_2 (0.43) - \frac{4}{7} \log_2 (0.57) \\&= -\frac{3}{7} (-1.21) - \frac{4}{7} (-0.81) \\&= 0.52 + 0.46 \\&= 0.97\end{aligned}$$

E<sub>normal</sub> = For Normal:

$$P = 6, N = 1$$

$$\begin{aligned}E_{\text{normal}} &= -\frac{6}{7} \log_2 \left(\frac{6}{7}\right) - \frac{1}{7} \log_2 \left(\frac{1}{7}\right) \\&= 0.19 + 0.40 \\&= 0.59\end{aligned}$$

Entropy of Humidity  $\propto$

$$\begin{aligned}E(\text{Humidity}) &= \frac{7}{14} (0.97) + \frac{7}{14} (0.59) \\&= 0.48 + 0.29 \\&= 0.77\end{aligned}$$

$$\text{Gain}(S, \text{humidity}) = I(P, n) - E(\text{humidity})$$

$$\begin{aligned}&= 0.940 - 0.77 \\&= 0.17\end{aligned}$$

Information gain for wind ( $\overset{①}{V^2}$ ,  $\overset{④}{V^-3}$ )

$$= E(P, T) - \sum_{v \in \{w, str\}} \frac{1/2 \cdot 1/2}{1/4} E(S_v)$$

$$= 0.94 - \left[ \frac{8}{14} E(6, 2) + \frac{6}{14} E(3, 3) \right]$$

$$= 0.94 - \left[ \frac{8}{14} \left( -\frac{6}{8} \log_2 \frac{6}{8} - \frac{2}{8} \log_2 \frac{2}{8} \right) \right]$$

$$\underline{\text{in}(S, \text{wind})} = \underline{0.048}$$

Information gain for temperature ( $\overset{②}{hot}$ ,  $\overset{③}{mild}$ ,  $\overset{④}{cool}$ )

$$= 0.94 - \left[ \frac{4}{14} E(2, 2) + \frac{6}{14} E(4, 2) + \frac{4}{14} E(3, 1) \right]$$

$$= 0.94 - \left[ \frac{4}{14} \left( -\frac{2}{4} \log_2 \left( \frac{2}{4} \right) - \frac{2}{4} \log_2 \left( \frac{2}{4} \right) \right) \right]$$

$$+ \frac{6}{14} \left( -\frac{4}{6} \log_2 \left( \frac{4}{6} \right) - \frac{2}{6} \log_2 \left( \frac{2}{6} \right) \right)$$

$$+ \frac{4}{14} \left( -\frac{3}{4} \log_2 \left( \frac{3}{4} \right) - \frac{1}{4} \log_2 \left( \frac{1}{4} \right) \right)$$

$$= 0.94 - 0.906$$

$$= \underline{0.029}$$

$$\text{Gain}(s, \text{outlook}) = 0.246$$

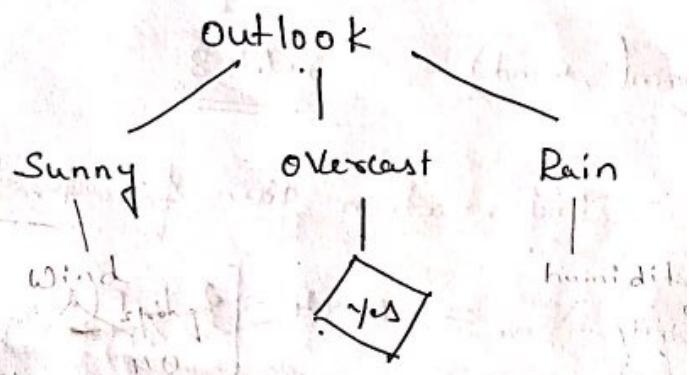
$$\text{Gain}(s, \text{humidity}) = 0.151$$

$$\text{Gain}(s, \text{wind}) = 0.048$$

$$\text{Gain}(s, \text{temp}) = 0.029$$

Comparing the info. gain for all the attributes, outlook attribute provides the best prediction of the target attribute.

Here outlook being overcast is always ideal yes for the target function.



Note:- Leaves of the decision tree represent class label

Branches represent conjunctions

Top most decision node in a tree corresponds to the best predictor called root node.

Apply recursion again with unadd attributes  
humidity ; temperature and rain left.

- \* Overcast is (value of outlook) is a yes, so now need to address sunny & rain.
- \* computing entropy for sunny.

Table for values of outlook & = sunny:

Temperature	Humidity	Wind	play tennis
Hot	high	weak	no
Hot	high	strong	no
Mild	high	weak	no
Cool	normal	weak	yes
mild	normal	strong	yes

$$E_{\text{sunny}} = -\frac{2}{5} \log_2 \left(\frac{2}{5}\right) - \frac{3}{5} \log_2 \left(\frac{3}{5}\right)$$

$$= 0.96$$

From previous steps compute information gain for humidity, temp and wind.

$$\text{Gain}(\text{sunny}, \text{humidity}) = 0.96$$

$$\text{Gain}(\text{sunny}, \text{Temp}) = 0.57$$

$$\text{Gain}(\text{sunny}, \text{wind}) = 0.019$$

Gain of sunny with humidity is high

Now calculate for rain

Final decision tree

