

What is data mining?

- nontrivial extraction of implicit, previously unknown and potentially useful information from data.
- exploration and analysis, by automatic or semi-automatic means, of large quantities of data in order to discover meaningful patterns.

What is not data mining?

- deductive query processing, we want is statistical algorithm running over the data.

Data mining application

1. Database analysis and decision support

↳ how market is segmented

↳ how the customers fill in the profile you address with the product.

↳ determine purchasing patterns.

↳ how different products are related.

2. Market analysis

↳ what the trends are.

↳ monitor competitors and market directions.

Input Data → Data Preprocessing → Data mining → Postprocessing → Information

• flat files, spreadsheet, relational tables	• raw input data to appropriate format for subsequent analysis	• closing the loop. • Integrating data mining results with decision support system.	• only valid and useful results are incorporated • explore from different view • filtering patterns, visualization
• centralized or distributed	• data cleaning • feature selection, dimensionality reduction, normalization, data subsetting		

Methods of Data mining

1. Predictive

↳ predict values of an attribute based on values of other attributes

↳ attribute to be predicted is known as target, while the other one is explanatory or independent.

2. Descriptive

↳ derive patterns that summarize relationships in data.

1. Predictive Modelling

→ building a model for the target variable as a function of independent variable.

(a) classification

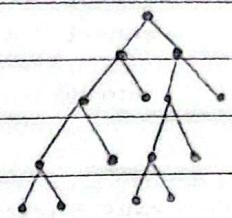
→ discrete target variables

→ predicting whether a web-user will make a purchase online is binary-valued.

(b) regression

→ continuous target variable.

→ forecasting future price of stock because price is continuous valued.



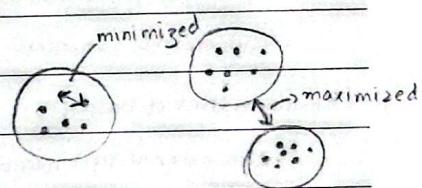
2. Cluster Analysis

→ groups of closely related observations so that observations that belong to the same cluster

are more similar to each other than observations of other cluster.

→ document clustering, similar articles grouped on a topic

→ K-means



3. Association Analysis

→ discover patterns that describe strongly associated features in the data.

→ represented in form of implication rules or feature subsets.

→ market-basket analysis

4. Anomaly Detection

→ identifying observations whose characteristics are significantly different from the rest of the data.

→ credit card fraud detection.

→ good anomaly detector must have high detection rate and a low false alarm rate.

What is data?

- collection of data objects and their attributes

- an attribute is a property or characteristic of an object

- attribute values are numbers or symbols assigned to an attribute

- height is attribute, feet or meters are values.

Types of Attributes

- nominal : used to categorise/classify objects, groups, individuals. e.g. gender, country, religion.
- ordinal : data is categorized and ranked, magnitude of rank unknown e.g. grades
- interval : categorize, rank and infer equal intervals between neighbouring data points e.g. temperature, calendar dates → mean, SD
- ratio : categorize, rank, infer equal intervals and there is a true zero point. e.g. weight, income, age, length, mass.
- discrete : has only finite set of values, integers, binary attributes e.g. ID number, zipcode.
- continuous : has real numbers as attribute value, floating point variables e.g. temperature, height
- symmetric : equal treatment of both values e.g. has a pet → Yes or No
- asymmetric : unequal treatment e.g. presence, non-zero attribute value is important

Characteristics of Data

- dimensionality : number of attributes.
- sparsity : those with asymmetric feature, attributes of an object have values of 0, presence of 1 counts
- resolution : low resolution data is historical gathered on hourly, monthly, yearly basis
high resolution data is collected at intervals of one-minute
- size : type of analysis depends on size.

Types of Dataset

1. Record
2. Graph
3. Ordered.

1. Record

- collection of records, each of which consists a fixed set of attributes.
- data matrix : each dimension represents a distinct attribute. m objects, n columns for each attribute.
- document data : each document becomes a 'term' vector, each term is a attribute and values is number of times the corresponding term occurs, non-zero entries are stored.
- transaction data : each transaction has set of items.

2. Graph

- relationship among data objects are links, and data objects are nodes.

3. Ordered Data

- sequential , each item purchased along with time.
- spatio temporal data.

Data Quality

- detection and correction of data quality problems
- data cleaning
- problems
 - ↳ noise, outliers, fake data, wrong data, missing values

1. Noise

- refers to modification of original values.
- techniques used to reduce noise to discover patterns/signal

2. Outliers

- data objects that are considerably different than most of the other data.
- anomalies
- outliers may sometimes be the goal of our analysis. e.g. credit card fraud.

Data Preprocessing

1. Aggregation

- combining two or more objects in a single object to reduce number of attributes.
- cities aggregated into regions, days to weeks or months.
- quantitative attributes e.g. price aggregated by sum, mean.
- have less variability.

2. Sampling

- selecting a subset of data objects to be analyzed.
- this is done because 'processing' entire set can be expensive.
- a sample is representative if it has approximately same properties as the original set of data.

(a) simple random

- equal probability of selecting any particular item.
- without replacement: as each item is selected, it is removed from population.
- with replacement: objects are not removed after selecting.

(b) stratified

- split data into several partitions and draw randomly.



KAGHAZ
www.kaghaz.pk

3. Discretization

- converts a continuous attribute into ordinal
- infinite number of values are mapped into small number of categories
- used in supervised and unsupervised.

4. Binarization

- converts discrete attribute into 1 or more binary attribute.

5. Attribute Transformation

- maps the entire set of values to a new set
- normalization
- standardization.

6. Dimensionality Reduction

- avoid curse of dimensionality
- reduce time, memory
- helps in eliminating irrelevant feature.
- use of PCA : creating new attributes that are combination of old one.
- captures largest amount of variation.

7. Feature Subset Selection

- redundant features duplicate much or all of the information - e.g.: price of product, sales tax paid.
- irrelevant features : information is not useful for data mining.

8. Feature Creation

- create new attributes that capture important information more efficiently than originals
- feature extraction : new features from raw data.
- feature construction.
- mapping data to new space

original value -

DM functionalities

• Association

→ one thing dependent on other.

→ $\text{age}(x, "20..29")$, $\text{income}(x, "20..29k") \rightarrow \text{buys}(x, "PC")$

People aged 20-29 with income 20-29k tends to buy PCs.

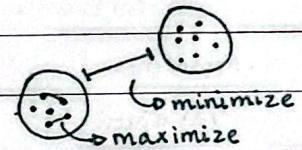
- Classification and Prediction

→ developing models like decision trees or neural networks for future predictions

• Cluster Analysis

→ class label is unknown

→ maximize intra-class similarity and minimize interclass similarity



• Outlier Analysis

→ a data object that is odd one out

→ detecting rare or fraudulent patterns in data.

Association Rule Mining

• finding all co-occurrence relationships among data items

• market basket analysis → discover how items are purchased.

$X \rightarrow Y$ (if X is bought, then Y is likely to be bought)

• Let $I = \{i_1, i_2, \dots, i_m\}$ be a set of items.

• Let $T = \{t_1, t_2, \dots, t_n\}$ be a set of transactions

$X \rightarrow Y$, where $X \subset I$, $Y \subset I$ and $X \cap Y = \emptyset$

someone who bought item X, also buys item Y e.g. Bread \rightarrow Butter but not Bread \rightarrow Bread

• Rules can be weak or strong

• Strength of rule is measured by its support and confidence

• Support \rightarrow probability that X and Y appear together in transaction.

• Confidence \rightarrow probability of Y appearing given X is present.

• Lift \rightarrow dependency between X and Y

$$\text{Lift}(X \Rightarrow Y) = \text{conf}(X \Rightarrow Y) / \text{supp}(Y)$$

Lift > 1 : strong association, appear together more than expected.

- if support is too low \rightarrow rules may appear by chance
- if confidence is low \rightarrow X does not reliably predict Y
- discover all associated rules in T that have support and confidence greater than a minimum threshold (minsup, minconf)

Approaches
 $\text{minsup} = 30\%$, $\text{minconf} = 80\%$.

$\text{chicken, clothes} \rightarrow \text{milk} = \frac{\text{chicken} \cap \text{clothes} \cap \text{milk}}{\text{no. of transaction}}$

$\text{sup} = \frac{3}{7}$

$\text{confidence} = \frac{\text{chicken} \cap \text{clothes} \cap \text{milk}}{\text{chicken} \cap \text{clothes}} = \frac{3}{3} = 1$

Transactions

T1 : beef, chicken, milk
T2 : beef, cheese
T3 : cheese, boots
T4 : beef, chicken, cheese
T5 : beef, chicken, clothes, milk, cheese
T6 : clothes, chicken, milk
T7 : chicken, milk, clothes.

Approaches

- apriori
- multiple minimum supports
- mining class association rule

Apriori Algorithm

Step 1: find all frequent itemsets (support > minsup)

Step 2: use frequent itemsets to generate candidate rules

STEP 1:	TRANSACTIONS	1. find frequent itemsets of size 1 $\rightarrow F_1$
T1	beef, chicken, milk	
T2	beef, cheese	
T3	cheese, boots	
T4	beef, chicken, cheese	
T5	beef, chicken, clothes, cheese, milk	
T6	clothes, chicken, milk	
T7	chicken, milk, clothes	

2. generalization, $K \geq 2$

C_K = those itemsets of size K that could be frequent.

F_K = those itemsets that are actually frequent

- Downward closure property \rightarrow any subset of a frequent itemset is also a frequent itemset.

example \rightarrow

$$\{\text{chicken, clothes, Milk}\} = \frac{3}{7}$$

$$(m) \text{ minsup} = 30\%$$

$$\frac{3}{7} = (\text{chicken, clothes}) \quad (\text{chicken, Milk}) = \frac{4}{7} \quad (\text{clothes, Milk}) = \frac{4}{7}$$

$$\frac{5}{7}, \frac{3}{7}, \frac{4}{7}$$

- Items are sorted in lexicographic order

• Generalization

step 1: joining (A_{k-1}, B_{k-1})

example $\rightarrow F_3 = \{\{1, 2, 3\}, \{1, 2, 4\}, \{1, 3, 4\}, \{1, 3, 5\}, \{2, 3, 4\}\}$

$1, 2, 3$	\rightarrow	$1, 2, 4$	\rightarrow	$1, 2, 3, 4$		$1, 2, 4$	\rightarrow	$1, 3, 4$
		$1, 3, 4$						$1, 3, 5$
		$1, 3, 5$						
		$2, 3, 4$						$2, 3, 4$

$1, 3, 4$	\rightarrow	$1, 3, 5$	\rightarrow	$1, 3, 4, 5$		$1, 3, 5$	\rightarrow	$2, 3, 4$
		$2, 3, 4$						

• only the last element should be different and greater than original set.

• in the above example there are 2 new candidates which could be frequent

step 2: pruning \rightarrow remove those candidates in C_k that do not respect the downward closure

property

$1, 2, 3, 4$	$1, 2, 3 \checkmark$	$1, 2, 4 \checkmark$	$1, 3, 4, 5$	$1, 3, 4 \checkmark$	$1, 3, 5 \checkmark$
	$1, 3, 4 \checkmark$			$1, 4, 5 X$	
	$2, 3, 4$			$3, 4, 5 X$	

$$C_4 = \{(1, 2, 3, 4), (1, 3, 4, 5)\}$$

all of these exist

remove from C_4

in F_{k-1} where $k=4$
so downward closure
property exist.

after pruning $C_4 = [1, 2, 3, 4]$

Multiple Minimum Support

- each item has its own minimum support
- some items appear frequently in the data, while others are rare.
- if minsup is too high, no rare items will be found
- to prevent very frequent items and very rare items from appearing in the same itemsets, we introduce a support difference constraint (α)

$$\max \{ \text{sup}(i) - \min \{ \text{sup}(i) \} \} \leq \phi \quad 0 \leq \phi \leq 1 \text{ is user specified.}$$

- MIS(i) is the minimum item support
- The minsup of a rule R is the lowest MIS value of all items.

Rule R satisfies its minimum support if;

$$\text{actual support} \geq \min(\text{MIS}(i_1), \text{MIS}(i_2), \dots)$$

- downward closure property doesn't exist anymore.

example \rightarrow MIS(1) = 10 · 1. MIS(2) = 20 · 1. MIS(3) = 5 · 1. MIS(4) = 6 · 1.

$\{1, 2\}$ has a support of 9 · 1. $\rightarrow 9 \geq \min(10 \cdot 1, 20 \cdot 1) \rightarrow$ not frequent

$\{1, 2, 3\}$ has a support of 7 · 1. $\rightarrow 7 \geq \min(10, 20, 5) \rightarrow$ frequent

* applying downward closure would eliminate $(1, 2)$ without evaluating $(1, 2, 3)$

which is frequent

→ How to solve downward closure property?

1. sort all items in I according to their MIS values, starting with smallest

→ Algorithm

Step 1: frequent itemset generation

(a) seeds for generating candidate itemsets.

(b) candidate generation for $k=2$

(c) generalization for $k > 2$

Step 2: rule generation



KAGHAZ
www.kaghaz.pk

STEP 1: frequent itemset generation

example $\rightarrow I = [1, 2, 3, 4]$ MIS(1) = 10%. MIS(2) = 20%. MIS(3) = 5%. MIS(4) = 6%. $n = 100$ transaction

(a). sort according to MIS $\rightarrow [3, 4, 1, 2]$

count $\rightarrow [3]: 6, [4]: 3, [1]: 9, [2]: 25$

- seeds of list L: {3}

$$\rightarrow \text{MIS}(3) = 5\%.$$

$$\text{item } 3 \rightarrow 6\% > 5\% \rightarrow \text{so } L = [3]$$

$$\text{item } 4 \rightarrow 3\% > 5\% \rightarrow \text{so } L = [3]$$

$$\text{item } 1 \rightarrow 9\% > 5\% \rightarrow \text{so } L = [3, 1]$$

$$\text{item } 2 \rightarrow 25\% > 5\% \rightarrow \text{so } L = [3, 1, 2]$$

. calculate F1 now if any item in L has actual support $< \text{MIS}(i)$

$$\text{item } 3 \rightarrow 6\% < 5\% \checkmark$$

$$\text{item } 1 \rightarrow 9\% < 10\% \times$$

$$\text{item } 2 \rightarrow 25\% < 20\% \checkmark$$

. so frequent itemset $= F_1 = [3, 2]$

(b) candidate generation, $k=2$, $\Phi = 10\%$.

* use L and not F₁ due to downward closure property.

. test chosen item against its MIS: $\text{sup}(3) > \text{MIS}(3)$

(a) if true, form level 2 candidates, otherwise go to next element

$$L = [3, 1, 2] \rightarrow [3, 1], [3, 2]$$

[3, 1] is a candidate $\rightarrow \text{sup}(1) > \text{MIS}(3)$ AND $|\text{sup}(3) - \text{sup}(1)| \leq \Phi$

$$9\% > 5\% \text{ AND } |6 - 9| \leq 10\%, \text{ thus } C_2 = [3, 1]$$

[3, 2] $\rightarrow \text{sup}(2) > \text{MIS}(3)$ AND $|\text{sup}(3) - \text{sup}(2)| \leq \Phi$

$$25\% > 6\% \text{ AND } |6 - 25| \leq 10\%, \text{ thus not true.}$$

. test: $\text{sup}(1) > \text{MIS}(1) = 9\% > 10\% = \text{not true so cannot proceed further. so } C_2 = [3, 1]$

. calculate support of each item in C₂ $\rightarrow \text{sup}(3, 1) = 6\% > \min[5, 10]$ so F₂ = [3, 1]



Mining Sequence Patterns

- frequently occurring ordered events
- customer retention, targeted marketing.
- an event is an itemset, unordered list of items, bought in a single transaction.

(I₁, I₂, I₃)

- sequence S is an ordered list of events

<e₁, e₂, e₃, ...>

- example:

e₁ = (abc), e₂ = (ade)S = <e₁, e₂> = <(abc)(ade)> → length of S = 6

- α is a subsequence of β because:

$$\alpha = \langle(ab)d\rangle \text{ and } \beta = \langle(abc)(de)\rangle \quad \therefore \alpha \sqsubseteq \beta$$

- support of a sequence

↳ number of tuples in S, containing α

- α is a frequent sequence if $\text{sup}_S(\alpha) \geq \text{min-sup}$.

- a sequence pattern is also a frequent sequence.

EXAMPLE

SID	SEQUENCE	I = [a,b,c,d,e,f,g]
1	<a(abc)(ac)a(cf)>	min-sup = 2
2	<(ad)c(bc)(ae)>	
3	<(ef)(ab)(df)cb>	
4	<eg(af)cbc>	

∴ for SID = 1, length = 9, but support of a remains 1 in first transaction even though it has occurred thrice, so total support = 4

∴ <a(bc)af> is a subsequence of SID = 1

$\langle a(bc)df \rangle$
 $\langle a(abc)(ac)d(cf) \rangle$

- huge number of sequential patterns are hidden in databases.



KAGHAZ
www.kaghazpk

- a priori based methods.

→ generalized sequential pattern. (based on downward closure property)

Generalized sequential pattern.

step 1: calculate support of all length 1 sequence and discard those which are less than min sup.

step 2: generate length 2 candidates from length 1 pattern. (ordered, unordered table)
1 event, 2 event

drawbacks

↳ huge set of candidate sequence generated

↳ multiple scan of database

↳ long patterns grow from short patterns.

- combining sequences of events with repeated measurements, we obtain time series data.

Time - Series Data

- reveal temporal behaviour

- events changing with time.

- data recorded at regular intervals.

- financial, industry, meteorological

- goals

→ modeling time-series

- analyze the past data.

→ forecasting time-series

- use the model to predict future data.

- methods

→ trend analysis

→ similarity search.

Trend Analysis

- statistical techniques. e.g. regression analysis.

- construct a model to explain the behavior of the measurement. e.g. correlation.

- regression analysis

↳ finding trends and outliers in datasets.

↳ dependent and independent variable.

- example

→ determine appropriate levels of advertising for a particular market segment

• linear regression analysis

$$Y = b_0 + b_1 X_1 + b_2 X_2 + \dots + b_n X_n$$

↑ difference on sales.

 Y = predicted score b_0 = intercept (origin) b_i = regression coefficients• correlation (R)

→ denoted by pearson

→ refers to interdependence.

→ closeness of linear relationship between X and Y .→ between -1 and 1
(positive)

error = $y - \hat{y}$

$\hat{y} = B_0 + B_1 X_1$

$MSE = \frac{1}{n} \sum (y_i - \hat{y}_i)^2$ actual - predicted

$B_0 = \bar{Y} - B_1 \bar{X}$

• regression trend channels

→ upper and lower trendline from standard deviation of regression line.

→ when the prices break a well established trend channel.

→ represented by 3 lines, where middle one is linear regression line and then \pm standard deviation.

• non-linear regression

→ bayesian methods

• characteristics time-series movement components .

1. trend → long term progression, upward downward

2. seasonal → identical patterns that a time series appears to follow. (christmas sales)

3. cycle → regular fluctuations (business cycle)

4. irregular → random

• decomposition

→ additive model = $T + C + S + I$ → multiplicative = $T \times C \times S \times I$ 

Trend Analysis (T), methods

1. freehand method

→ fit the curve by looking at the graph.

→ costly

2. least-square method

→ curve minimizing the sum of the squares.

3. moving average method

→ eliminates cyclic, seasonal and irregular patterns.

→ loss of ends as it is a sliding window.

→ sensitive outliers.

→ a smoother curve

→ influence of extreme values can be reduced with weighted moving average (WMA)

Original WMA(3) → (141)

3

7

2

0

4

5

9

7

2

$$(3 \times 1 + 4 \times 7 + 1 \times 2) / 1+4+1$$

$$(7 \times 1 + 2 \times 4 + 0 \times 1) / 1+4+1$$

→ cumulative moving average

$$CA_i = \frac{x_1 + \dots + x_i}{i}, \quad CA_{i+1} = CA_i + \frac{x_{i+1} - CA_i}{i+1}$$

→ exponential weighted moving average.

Estimation of seasonal variation (s)

1. seasonal index

* happens regularly at fixed points within a year

2. deseasonalized

→ adjusted with seasonal variations.

Estimation of cyclic variation (c)

→ semantic meaning

* happens over irregular time frames, tied to long term economic factor.

Estimation of irregular variation (e)

→ adjusting the data.

• make long or short term predictions (time-series forecasting)



KAGHAZ
www.kaghaz.pk

day / date:

- time series forecasting

- Auto regressive integrated moving average

- non-stationary.

- granularity change

- moving average for trend analysis

- ↳ identify open position for buying and selling

- bollinger bands → change of trends.

- resistance lines → market stops at one point

- momentum analysis.

- ↳ difference between opening and closing position.

Similarity Search

- exact matches in normal database

- SS finds data sequence that differ only slightly

- financial market, market basket

Classification

- collection of records
- each record consists of a set of attributes
 - x : attribute, independent, input
 - y : class, response, dependent, output
- goal is to assign new records to each class without asking
- defining a class attribute as a function of characteristics that you know
- training on a known data and then testing to validate
- financial field, news
- training set → learn model → apply model → deduce → test set → classification.

unknown class

Supervised learning

- training data comes with class labels
- new data classified based on training set

Unsupervised learning

- no labels in training set
- clustering

Decision Trees : Base classifiers

internal node : test on an attribute

branch : represents outcome of test

leafnodes : class labels

STEP 1 : INDUCTION

- input : training dataset, attribute list
- output : decision tree.

STEP 2 : DEDUCTION

- predict classes of new data
- input : decision tree, new data
- output : classified new data

DECISION TREE INDUCTION

1. HUNT'S ALGORITHM

- $D_t \rightarrow$ set of training records (data) that reach a node t
- split the data in classes on an attribute test
 - ↳ nominal (discrete values)
 - ↳ ordinal (categories)
 - ↳ continuous (numbers) \rightarrow range, buckets.
 - ↳ binary split
 - ↳ multi-way split
- if D_t contains records that belong to same class y_t then t is a leafnode labelled as y_t .
- If D_t contains records that belong to more than one class, choose next best attribute for test condition and recursively apply it.
- If D_t is not an empty set, but no other attributes remain for splitting then t is a leaf node labelled by label of majority records.
- splitting continuous variable
 - ↳ discretization (static, dynamic)
 - ↳ binary decision (interval, right or left)
- to determine the best split
 - ↳ check ratios where it is greatest
 - ↳ nodes with homogenous class distribution preferred

example: student \rightarrow Yes: Yes: 6, No: 1 ✓
 income \rightarrow high, Yes: 2, No: 2 ✗

Methods to measure impurity

1. Information Gain

ratio of classes of some label opposed to classes of other labels

Yes: 6 No: 1 \rightarrow high

Yes: 2 No: 2 \rightarrow low

2. Gini Index

• continuous valued, several splits

1. Information Gain

assume there are 2 classes, P and N

p elements of class P, n elements of class N

$$I(p,n) = \frac{-p}{p+n} \log_2 \frac{p}{p+n} - \frac{n}{p+n} \log_2 \frac{n}{p+n}$$

probability that label is p probability that label is N

information gain in decision tree

↪ partition into sets.

$$E(A) = \sum_{i=1}^k \frac{p_i + n_i}{p+n} I(p_i, n_i)$$

sum of entropy of all subsets

↪ encoding information

$$\text{Gain}(A) = I(p,n) - E(A)$$

↪ bigger entropy, worst split

EXAMPLE - STEP 1

class P : buys computer = Yes = 9

class N : buys computer = No = 5

1. $I(p,n) = I(9,5) = \frac{-9}{9+5} \log_2 \frac{9}{9+5} - \frac{5}{9+5} \log_2 \frac{5}{9+5} = 0.94$

2. compute entropy on basis of age

age	p_i	n_i	$I(p_i, n_i)$
≤ 30	2	3	0.971
> 40	4	0	0
$31 \dots 40$	3	2	0.971

$$= \sum \frac{2+3}{14} I(2,3) + \frac{4+0}{14} I(4,0) + \frac{3+2}{14} I(3,2) = 0.694$$

3. calculate $\text{gain}(\text{age}) = I(9,5) - E(\text{age}) = 0.94 - 0.694 = 0.246$

→ higher important, pure group, less entropy, more information gain.

$\text{gain}(\text{age}) > \text{income, student, credit}$

4. now remove age and make further splits on income, student, credit.

keep which has higher gain.

STOP CONDITION

- When all records belong to same class
- Or when similar attribute values, a leaf node labelled as majority class

STEP 2: DEDUCTION

- Read from root node and reach leafnode -
- One rule created for each path from root to leaf .
- attribute - value pair forms a conjunction .

Advantages

- inexpensive to construct
- fast in classifying unknown records.
- easy to interpret for small sized trees .
- avoid overfitting in classification .
 - ↳ prepruning → do not split if IG falling below a threshold
 - ↳ post pruning → remove branches from fully grown trees .

Mining Sequence Patterns

- frequently occurring ordered events
- customer retention, targeted marketing.
- an event is an itemset, unordered list of items, bought in a single transaction.

 (I_1, I_2, I_3)

- sequence S is an ordered list of events

 $\langle e_1, e_2, e_3, \dots \rangle$

- example;

$$e_1 = (abc), e_2 = (ade)$$

$$S = \langle e_1, e_2 \rangle = \langle (abc)(ade) \rangle \rightarrow \text{length of } S = 6$$

- α is a subsequence of β because;

$$\alpha = \langle (ab)d \rangle \text{ and } \beta = \langle (abc)(de) \rangle$$
$$\therefore \alpha \sqsubseteq \beta$$

- support of a sequence

\hookrightarrow number of tuples in S, containing α

- α is a frequent sequence if $\text{sup}_{\text{S}}(\alpha) \geq \text{min-sup}$

- a sequence pattern is also a frequent sequence.

EXAMPLE

SID	SEQUENCE	$I = [a, b, c, d, e, f, g]$
1	$\langle a(abc)(ac)d(cf) \rangle$	$\text{min-sup} = 2$
2	$\langle (ad)c(bc)(ae) \rangle$	
3	$\langle (ef)(ab)(df)cb \rangle$	
4	$\langle eg(af)cbc \rangle$	

\therefore for SID = 1, length = 9, but support of α remains 1 in first transaction even though it has occurred thrice, so total support = 4

$\therefore \langle a(bc)af \rangle$ is a subsequence of SID = 1

$\langle a(bc)af \rangle \sqsubseteq \langle a(abc)(ac)d(cf) \rangle$

- huge number of sequential patterns are hidden in databases.

- a priori based methods

→ generalized sequential pattern (based on downward closure property)

Generalized sequential pattern.

step 1: calculate support of all length 1 sequence and discard those which are less than minsup.

step 2: generate length 2 candidates from length 1 pattern (ordered, unordered table)
1 event, 2 event

drawbacks

↳ huge set of candidate sequence generated

↳ multiple scan of database

↳ long patterns grow from short patterns.

• combining sequences of events with repeated measurements, we obtain time series data.

Time - Series Data

- reveal temporal behaviour

- events changing with time

- data recorded at regular intervals

- financial, industry, meteorological

- goals

→ modeling time-series

- analyze the past data.

→ forecasting time-series

- use the model to predict future data.

- methods

→ trend analysis

→ similarity search.

Trend Analysis

- statistical techniques. e.g. regression analysis.

- construct a model to explain the behavior of the measurement. e.g. correlation.

- regression analysis

↳ finding trends and outliers in datasets

↳ dependent and independent variable.

- example

→ determine appropriate levels of advertising for a particular market segment.

• linear regression analysis

$$Y = b_0 + b_1 X_1 + b_2 X_2 + \dots + b_n X_n$$

Y = predicted score

b_0 = intercept / origin

b_i = regression coefficients

• correlation

→ denoted by pearson

→ refers to interdependence.

→ closeness of linear relationship between X and Y .

→ between -1 and 1
(positive)

$$\text{error} = y - \hat{y}$$

$$\hat{y} = B_0 + B_1 X_1$$

$$\text{MSE} = \frac{1}{n} \sum (y_i - \hat{y}_i)^2$$

$$B_0 = \bar{y} - B_1 \bar{x}$$

• regression trend channels

→ upper and lower trendline from standard deviation of regression line.

→ when the prices break a well established trend channel.

→ represented by 3 lines, where middle one is linear regression line and then \pm standard deviation

• non-linear regression

→ bayesian methods

• characteristics time-series movement components

1. trend → long term progression

2. seasonal → identical patterns that a time series appears to follow. (christmas sales)

3. cycle → regular fluctuations (business cycle)

4. irregular → random

• decomposition

→ additive model = $T + C + S + I$

→ multiplicative = $T \times C \times S \times I$

Trend Analysis (T), methods

1. freehand method

→ fit the curve by looking at the graph.

→ costly

2. least-square method

→ curve minimizing the sum of the squares.

3. moving average method

→ eliminates cyclic, seasonal and irregular patterns

→ loss of ends as it is a sliding window.

→ sensitive outliers.

→ a smoother curve

→ influence of extreme values can be reduced with weighted moving average (WMA)

Original

WMA(3) → (141)

3

7

2

0

4

5

9

7

2

$$(3 \times 1 + 4 \times 7 + 1 \times 2) / 1+4+1$$

$$(7 \times 1 + 2 \times 4 + 0 \times 1) / 1+4+1$$

→ cumulative moving average

$$CA_i = \frac{x_1 + \dots + x_i}{i}, \quad CA_{i+1} = CA_i + \frac{x_{i+1} - CA_i}{i+1}$$

→ exponential weighted moving average.

Estimation of seasonal variation (s)

1. seasonal index

* happens regularly at fixed points within a year

2. deseasonalized

→ adjusted with seasonal variations.

Estimation of cyclic variation (c)

→ semantic meaning

* happens over irregular time frames, tied to long term economic factor.

Estimation of irregular variation (i)

→ adjusting the data.

- make long or short term predictions (time-series forecasting)



KAGHAZ
www.kaghaz.pk

- time series forecasting
 - Auto regressive integrated moving average.
 - non-stationary.
- granularity change
- moving average for trend analysis
 - ↳ identify open position for buying and selling.
- bollinger bands
- resistance lines
- momentum analysis.
 - ↳ difference between opening and closing position.

Similarity Search

- exact matches in normal database
- SS finds data sequence that differ only slightly
- financial market, market basket

Classification

- collection of records
- each record consists of a set of attributes
 - x : attribute, independent, input
 - y : class, response, dependent, output
- goal is to assign new records to each class without asking
- defining a class attribute as a function of characteristics that you know.
- training on a known data and then testing to validate.
- financial field, news
- training set → learn model → apply model → deduce → test set → classification.
unknown class

Supervised learning

- training data comes with class labels
- new data classified based on training set

Unsupervised learning

- no labels in training set
- clustering

Decision Trees : Base classifiers

internal node : test on an attribute

branch : represents outcome of test

leafnodes: class labels.

STEP 1: INDUCTION

- input: training dataset, attribute list
- output: decision tree.

STEP 2: DEDUCTION

- predict classes of new data.
- input: decision tree, new data
- output: classified new data.

DECISION TREE INDUCTION

1. HUNT'S ALGORITHM

$D_t \rightarrow$ set of training records (data) that reach a node t

- split the data in classes on an attribute test

↳ nominal (discrete values)

↳ ordinal (categories)

↳ continuous (numbers)

↳ binary split

↳ multi-way split

• if D_t contains records that belong to same class y , then t is a leafnode labelled as y_t .

• if D_t contains records that belong to more than one class, choose next best attribute for test condition and recursively apply it.

• if D_t is not an empty set, but no other attributes remain for splitting then t is a leafnode labelled by label of majority records

- splitting continuous variable

↳ discretization (static, dynamic)

↳ binary decision (interval, right or left)

- to determine the best split

↳ check variation where it is greatest

↳ nodes with homogenous class distribution preferred

example: student \rightarrow Yes: 6, No: 1 ✓

income \rightarrow high, Yes: 2, No: 2 ✗

Methods to measure impurity

1. Information Gain

ratio of classes of some label opposed to classes of other labels

Yes: 6 No: 1 \rightarrow high

Yes: 2 No: 2 \rightarrow low

2. Gini Index

- continuous valued, several splits

1. Information Gain

- assume there are 2 classes, P and N

- p elements of class P, n elements of class N

$$I(p,n) = \frac{-p}{p+n} \log_2 \frac{p}{p+n} - \frac{n}{p+n} \log_2 \frac{n}{p+n}$$

↓ ↓

probability
that label is P probability
that label is N

- information gain in decision tree

↳ partition into sets

$$E(A) = \sum_{i=1}^n \frac{p_i + n_i}{p+n} I(p_i, n_i)$$

sum of entropy of all subsets

↳ encoding information

$$\text{Gain}(A) = I(p,n) - E(A)$$

- bigger entropy, worst split

EXAMPLE - STEP 1

Class P : buys computer = Yes = 9

Class N : buys computer = No = 5

1. $I(p,n) = I(9,5) = \frac{-9}{9+5} \log_2 \frac{9}{9+5} - \frac{5}{9+5} \log_2 \frac{5}{9+5} = 0.94$

2. compute entropy on basis of age

age	p_i	n_i	$I(p_i, n_i)$
≤ 30	2	3	0.971
> 40	4	0	0
$31 \dots 40$	3	2	0.971

$$= \sum \frac{2+3}{14} I(2,3) + \frac{4+0}{14} I(4,0) + \frac{3+2}{14} I(3,2) = 0.694$$

3. calculate $\text{gain}(\text{age}) = I(9,5) - E(\text{age}) = 0.94 - 0.694 = 0.246$

→ higher important, pure group, less entropy, more information gain.

$\text{gain}(\text{age}) > \text{income, student, credit}$

4. now remove age and make further splits on income, student, credit.

Keep which has higher gain.

STOP CONDITION

- when all records belong to same class.
- or when similar attribute values, a leaf node labelled as majority class

STEP 2: DEDUCTION

- read from root node and reach leafnode.

- one rule created for each path from root to leaf.

- attribute - value pair forms a conjunction.

Advantages

- inexpensive to construct
- fast in classifying unknown records
- easy to interpret for small sized trees

→ avoid overfitting in classification.

↳ prepruning → do not split if IG falling below a threshold

↳ post pruning → remove branches from fully grown trees.

SEQUENCE PATTERN MINING

SID	SEQUENCE
1	< milk(bread,butter)(eggs)diapers >
2	<(milk,eggs)(bread,butter)diapers >
3	<(cereal,milk)juice(bread,butter)eggs >
4	<(eggs,juice)(bread,butter)diapers >
5	<(cereal,milk)eggs(juice,butter) >
6	<juice(milk,bread)(butter,eggs)diapers >
7	<(cereal,juice)(bread,eggs)(butter)milk >

STEP 1 : CREATE 1-LENGTH SEQUENCE

$\langle \text{milk} \rangle = a = 6$, $\langle \text{bread} \rangle = b = 6$, $\langle \text{butter} \rangle = c = 7$, $\langle \text{eggs} \rangle = d = 7$

$\langle \text{diapers} \rangle = e = 4$, $\langle \text{cereal} \rangle = f = 3$, $\langle \text{juice} \rangle = g = 5$

STEP 2 : 1-LENGTH PATTERN

$\langle a \rangle = 6$ $\langle b \rangle = 6$ $\langle c \rangle = 7$ $\langle d \rangle = 7$ $\langle e \rangle = 4$ $\langle g \rangle = 5$ $\langle f \rangle = 3$

• discard none as all over > min-sup.

STEP 3 : 2-LENGTH SEQUENCE

$\langle ab \rangle$	$\langle ba \rangle$	$\langle ac \rangle$	$\langle ca \rangle$	$\langle cd \rangle$	$\langle dc \rangle$	$\langle ce \rangle$	$\langle ec \rangle$	$\langle cf \rangle$	$\langle fc \rangle$	$\langle dg \rangle$	$\langle gd \rangle$
a^0	a^3	a^5	a^4	a^3	a^0	a^2	a^0	a^0	a^0	a^0	a^0
b^0	b^1	b^0	b^2	b^3	b^4	b^0	b^1	b^0	b^1	b^0	b^0
c^0	c^1	c^0	c^0	c^2	c^4	c^1	c^0	c^0	c^1	c^0	c^0
d^0	d^1	d^2	d^3	d^4	d^0	d^1	d^0	d^1	d^0	d^1	d^0
e^0	e^1	e^0	e^1	e^2	e^3	e^0	e^1	e^0	e^1	e^0	e^1
f^0	f^1	f^2	f^3	f^4	f^0	f^1	f^0	f^1	f^0	f^1	f^0
g^0	g^1	g^2	g^3	g^4	g^0	g^1	g^0	g^1	g^0	g^1	g^0

$\langle ab \rangle$	$\langle ba \rangle$	$\langle ac \rangle$	$\langle ca \rangle$	$\langle cd \rangle$	$\langle dc \rangle$	$\langle ce \rangle$	$\langle ec \rangle$	$\langle cf \rangle$	$\langle fc \rangle$	$\langle dg \rangle$	$\langle gd \rangle$
$(ab)^0$	$(ab)^1$	$(ac)^0$	$(ac)^1$	$(ad)^0$	$(ad)^1$	$(ae)^0$	$(ae)^1$	$(af)^0$	$(af)^1$	$(ag)^0$	$(ag)^1$
$(ba)^0$	$(ba)^1$	$(bc)^0$	$(bc)^1$	$(bd)^0$	$(bd)^1$	$(be)^0$	$(be)^1$	$(bf)^0$	$(bf)^1$	$(bg)^0$	$(bg)^1$
$(ca)^0$	$(ca)^1$	$(cb)^0$	$(cb)^1$	$(cd)^0$	$(cd)^1$	$(ce)^0$	$(ce)^1$	$(cf)^0$	$(cf)^1$	$(cg)^0$	$(cg)^1$
$(da)^0$	$(da)^1$	$(db)^0$	$(db)^1$	$(dc)^0$	$(dc)^1$	$(de)^0$	$(de)^1$	$(df)^0$	$(df)^1$	$(dg)^0$	$(dg)^1$
$(ea)^0$	$(ea)^1$	$(eb)^0$	$(eb)^1$	$(ec)^0$	$(ec)^1$	$(ed)^0$	$(ed)^1$	$(ef)^0$	$(ef)^1$	$(eg)^0$	$(eg)^1$
$(fa)^0$	$(fa)^1$	$(fb)^0$	$(fb)^1$	$(fc)^0$	$(fc)^1$	$(fd)^0$	$(fd)^1$	$(fe)^0$	$(fe)^1$	$(fg)^0$	$(fg)^1$
$(ga)^0$	$(ga)^1$	$(gb)^0$	$(gb)^1$	$(gc)^0$	$(gc)^1$	$(gd)^0$	$(gd)^1$	$(ge)^0$	$(ge)^1$	$(gf)^0$	$(gf)^1$
$(gb)^0$	$(gb)^1$	$(gc)^0$	$(gc)^1$	$(gd)^0$	$(gd)^1$	$(ge)^0$	$(ge)^1$	$(gf)^0$	$(gf)^1$	$(gg)^0$	$(gg)^1$

$\langle ab \rangle$ $\langle ac \rangle$ $\langle ad \rangle$ $\langle ae \rangle$
 $\langle bd \rangle$ $\langle be \rangle$
 $\langle ce \rangle$
 $\langle dc \rangle$ $\langle de \rangle$
 $\langle fc \rangle$ $\langle fd \rangle$
 $\langle gb \rangle$ $\langle gc \rangle$ $\langle gd \rangle$

$\langle (bc) \rangle$
 $\langle (bf) \rangle$
 $\langle (ca) \rangle$
 $\langle (bc) \rangle$

1. max speedup

Naive Bayes

age	income	student	credit rating	buys comp.
≤ 30	H	No	fair	No
≤ 30	H	No	excellent	No
31..40	H	No	fair	Yes
> 40	M	No	fair	Yes
> 40	L	Yes	fair	Yes
> 40	L	Yes	excellent	No
31..40	L	Yes	excellent	Yes
≤ 30	M	No	fair	No
≤ 30	L	Yes	fair	Yes
> 40	M	Yes	fair	Yes
≤ 30	M	Yes	excellent	Yes
31..40	M	No	excellent	Yes
31..40	H	Yes	fair	Yes
> 40	M	No	excellent	No

1. calculate probability for target column.

$$\text{positive (P)} = \text{buys - computer} = \text{Yes} = 9/14$$

$$\text{negative (n)} = \text{buys - computer} = \text{No} = 4/14$$

2. calculate probability for "age" column,
for each interval corresponding to Yes and No
in target column.

age	yes	no
≤ 30	2/9	3/5
31..40	4/9	0/5
> 40	3/9	2/5

3. calculate for remaining columns:

income, student, credit rating

income	yes	no
H	2/9	2/5
M	4/9	2/5
L	3/9	1/5

student	yes	no
Yes	3/9	1/5
No	3/9	4/5

credit rating	Yes	No
fair	6/9	2/5
excellent	3/9	3/5

4. classify according to given instructions

x = Age: youth

Income: low

Student: yes

credit: fair

separate for Yes and No

$$[P(x|p) * P(p)] / P(x) \rightarrow \text{positive}$$

$$= \frac{2}{9} \times \frac{3}{9} \times \frac{6}{9} \times \frac{6}{9} \times \frac{9}{14}$$

$$= 0.02111$$

$$[P(x|n) * P(n)] / P(x) \rightarrow \text{negative}$$

$$= \frac{3}{5} \times \frac{1}{5} \times \frac{1}{5} \times \frac{2}{5} \times \frac{4}{14}$$

$$= 0.0034$$

5. conclude and classify

since value of positive > value of negative,
so classify it as "buys-computer".

ID	Items Purchased
1	Milk, Bread, Butter, Eggs
2	Milk, Bread, Butter
3	Milk, Eggs, Bread
4	Bread, Butter
5	Milk, Cereal, Bread
6	Jam, Eggs
7	Milk, Butter, Eggs
8	Milk, Bread, Butter
9	Jam, Cereal
10	Milk, Cereal

Item	Minsupport	Sup	minsupp =
(1) Milk	60%	7/10	minsupp = 50%
(2) Bread	50%	6/10	$\phi = 10\%$
(3) Eggs	50%	4/10	
(4) Butter	40%	5/10	
(5) Jam	25%	2/10	
(6) Cereal	10%	3/10	

sort according to MIS

[6, 5, 4, 3, 2, 1]

→ seeds of list L: {3 → MIS(6) = 10%}

[1:7], [2:6], [3:4], [4:5],

[5:2], [6:3]

→ check if $sup(j) \geq MIS(j)$

item 6 → $3 \geq 10\% \rightarrow$ so L = [6]

item 5 → $2 \geq 10\% \rightarrow$ so L = [6, 5]

item 4 → $5 \geq 10\% \rightarrow$ so L = [6, 5, 4]

item 3 → $4 \geq 10\% \rightarrow$ so L = [6, 5, 4, 3]

item 2 → $6 \geq 10\% \rightarrow$ so L = [6, 5, 4, 3, 2]

item 1 → $7 \geq 10\% \rightarrow$ so L = [6, 5, 4, 3, 2, 1]

→ calculating F₁ from L based on MIS of each item in L

item 6 → $3 \geq 10 \rightarrow \checkmark$

item 5 → $2 \geq 25 \rightarrow \times$

item 4 → $5 \geq 40 \rightarrow \checkmark$

item 3 → $4 \geq 50 \rightarrow \times$

item 2 → $6 \geq 50 \rightarrow \checkmark$

item 1 → $7 \geq 60 \rightarrow \checkmark$

F₁ = [6, 4, 2, 1]

→ continuing L1 [6, 5, 4, 3, 2, 1]

[6, 5] is a candidate → $sup(5) \geq MIS(6) \text{ AND } |sup(6) - sup(5)| \leq \phi$
 $\frac{2}{10} \geq \frac{10}{100} \text{ AND } |3-2| \leq \frac{10}{100}, \text{ true}$

[6, 4] is a candidate → $sup(4) \geq MIS(6) \text{ AND } |sup(6) - sup(4)| \leq \phi$
 $\frac{5}{10} \geq \frac{10}{100} \text{ AND } |3-5| \leq \frac{10}{100}, \text{ false}$

[6, 3] → $sup(3) \geq MIS(6) \text{ AND } |sup(6) - sup(3)| \leq \phi$
 $4 \geq 10 \text{ AND } |3-4| \leq \frac{10}{100}, \text{ true}$

[6, 2] → $6 \geq 10 \text{ AND } |3-6| \leq \frac{10}{100}, \text{ false}$

[6, 1] → $7 \geq 10 \text{ AND } |3-7| \leq \frac{10}{100}, \text{ false}$

sup(j)

discard 5

[4, 3] → $sup(3) \geq MIS(4) \text{ AND } |sup(4) - sup(3)| \leq \phi$
 $4 \geq 40 \text{ AND } |5-4| \leq 0.1, \text{ true}$

[4, 2] → $sup(2) \geq MIS(4) \text{ AND } |sup(4) - sup(2)| \leq \phi$
 $6 \geq 40 \text{ AND } |5-6| \leq 0.1, \text{ true}$

[4, 1] → $sup(1) \geq MIS(4) \text{ AND } |sup(4) - sup(1)| \leq \phi$
 $7 \geq 40 \text{ AND } |5-7| \leq 0.1, \text{ false}$

discard 3

[2, 1] → $sup(1) \geq MIS(2) \text{ AND } |sup(2) - sup(1)| \leq \phi$
 $7 \geq 50 \text{ AND } |6-7| \leq 0.1, \text{ true}$

C₂ = [6, 5], [6, 3], [4, 3], [4, 2], [2, 1]

→ calculate support

[6, 5]: 1, [6, 3]: 0, [4, 3]: 2, [4, 2]: 4

[2, 1]: 5

→ for F₂, check $sup(i, j) \geq \min(MIS(i), MIS(j))$

[6, 5] = 1 $\geq \min(10/25)$ valid

[6, 3] = 0 $\geq \min(10, 50)$ invalid

[4, 3] = 2 $\geq \min(40, 50)$ invalid

[4, 2] = 4 $\geq \min(40, 50)$ valid

[2, 1] = 5 $\geq \min(50, 60)$ valid

thus F₂ = [6, 5], [4, 2], [2, 1]

Q. 2)
 Q. 3)

ID	Items Purchased
1	(1) Milk, Bread, Butter ⁽²⁾
2	Bread, Butter, Eggs ⁽⁴⁾
3	Milk, Bread, Butter, Eggs
4	Milk, Bread
5	Bread, Butter
6	Milk, Bread, Butter, Eggs
7	Bread, Eggs
8	Milk, Eggs
9	Milk, Bread, Butter
10	Milk, Bread, Eggs

$[Milk, Eggs] \rightarrow [Bread, Butter] \rightarrow \text{none}$
 $[Bread, Eggs]$
 $[Butter, Eggs]$
 $[Bread, Butter] \rightarrow [Bread, Eggs] \rightarrow [Bread, Butter, Eggs]$
 $[Bread, Eggs] \rightarrow [Butter, Eggs]$

→ pruning

$C_3 = [Milk, Bread, Butter], [Milk, Bread, Eggs],$
 $[Milk, Butter, Eggs], [Bread, Butter, Eggs]$

all items in C_3 belong to F_2

$\xrightarrow{3} C_3 = [Milk, Bread, Butter : 4],$
 $[Milk, Bread, Eggs : 3]$
 $[Milk, Butter, Eggs : 2]$
 $[Bread, Butter, Eggs : 3]$

$F_3 = [Milk, Bread, Butter : 4],$
 $[Milk, Bread, Eggs : 3],$
 $[Bread, Butter, Eggs : 3]$

→ join

$C_4 = [Milk, Bread, Butter] \rightarrow [Milk, Bread, Eggs]$
 $[Bread, Butter, Eggs]$
 $[Milk, Bread, Eggs] \rightarrow [Bread, Butter, Eggs] \rightarrow \text{none}$

→ pruning

$C_4 = [Milk, Bread, Butter, Eggs]$
 all items in C_4 belong to F_3

$\xrightarrow{4} C_4 = [Milk, Bread, Butter, Eggs : 2]$

$F_4 = \text{none items.}$

$F_2 = \text{all items}$

→ join, make triples

$C_3 = [Milk, Bread] \rightarrow [Milk, Butter] \rightarrow [Milk, Bread, Butter]$
 $[Milk, Eggs]$
 $[Bread, Butter] \rightarrow [Milk, Bread, Eggs]$
 $[Bread, Eggs]$
 $[Butter, Eggs]$

$[Milk, Butter] \rightarrow [Milk, Eggs] \rightarrow [Milk, Butter, Eggs]$
 $[Bread, Butter]$
 $[Bread, Eggs]$
 $[Butter, Eggs]$

DECISION TREE

OUTLOOK	TEMP	HUMIDITY	WINDY	PLAY GOLF
Rainy	Hot	High	False	No
Rainy	Hot	H	True	No
Overcast	Hot	H	F	Yes
Sunny	Mild	H	F	Y
Sunny	Cool	Normal	F	Y
Sunny	Cool	N	T	N
Overcast	Cool	N	T	Y
Rainy	Mild	H	F	N
Rainy	Cool	N	F	Y
Sunny	Mild	N	F	Y
Rainy	Mild	N	T	Y
Overcast	Mild	H	T	Y
Overcast	Hot	N	F	Y
Sunny	Mild	H	T	N

1. calculate information gain for target class

$$\text{class P} = \text{play golf} = 9$$

$$\text{class N} = \text{does not} = 5$$

$$I(9,5) = \frac{-(9)}{14} \log_2 \frac{9}{14} - \frac{5}{14} \log_2 \frac{5}{14}$$

$$= 0.94$$

2. compute entropy for each column and check which column has highest gain, split on that column.

Pi	ni	I(pi, n)
R	2	3
O	4	0
S	3	2

$$= \sum \frac{2+3}{14} I(2,3) + \frac{4+0}{14} I(4,0) + \frac{3+2}{14} I(3,2)$$

$$= \frac{5}{14} (0.97) + \frac{4}{14} (0) + \frac{5}{14} (0.97) = 0.694$$

$$\text{gain(outlook)} = 0.94 - 0.694 = 0.246$$

Pi	ni	I(pi, n)
H	2	2
M	4	2
C	3	1

$$= \sum \frac{4}{14} (1) + \frac{6}{14} (0.918) + \frac{4}{14} (0.81) = 0.918$$

$$\text{gain(temp)} = 0.94 - 0.918 = 0.022$$

Pi	ni	I(pi, n)
H	3	4
N	6	1

$$= \sum \frac{7}{14} (0.985) + \frac{7}{14} (0.591) = 0.788$$

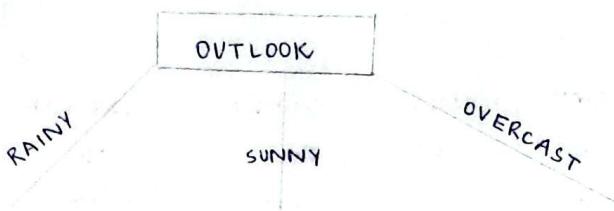
$$\text{gain(humidity)} = 0.94 - 0.788 = 0.152$$

Pi	ni	I(pi, n)
F	6	2
T	3	3

$$= \sum \frac{8}{14} (0.811) + \frac{6}{14} (1) = 0.891$$

$$\text{gain(windy)} = 0.94 - 0.891 = 0.048$$

3. outlook has highest gain so split on that.



T	H	W	
Hot	High	False	N
Hot	High	True	N
Mild	High	False	N
Cool	Norm	False	Y
Mild	Norm	True	Y

T	H	W	
Mild	High	False	Y
Cool	Norm	False	Y
Cool	Norm	True	N
Mild	Norm	False	Y
Mild	High	True	N

T	H	W	
Hot	High	False	Y
Cool	Norm	True	Y
Mild	High	True	Y
Hot	Norm	False	Y

$$I(2,3) = 0.97$$

$$I(3,2) = 0.971$$

✓ play-golf = yes.

for rainy =>

	Pi	ni	I(pi,n)
H	0	2	0
M	1	1	1
C	1	0	0

$$= \sum \frac{2}{5} (1) = 0.4$$

$$\text{gain(temp)} = 0.97 - 0.94 = 0.03$$

	Pi	ni	I(pi,n)
H	0	3	0
N	2	0	0

✓ gain(temp) = 0.97 - 0 = 0.97

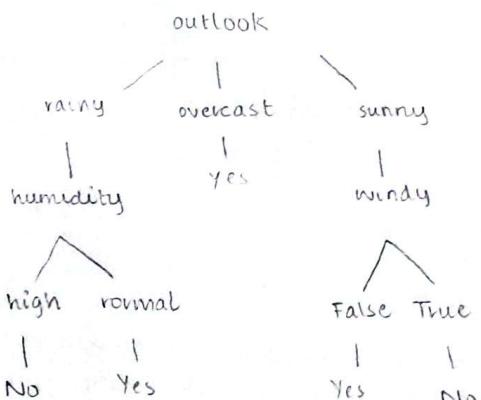
	Pi	ni	I(pi,n)
F	1	2	0.918
T	1	1	1

$$= \sum \frac{3}{5} (0.918) + \frac{2}{5} (1) = 0.9508$$

$$\text{gain(windy)} = 0.97 - 0.9508 = 0.0192$$

	Pi	ni	I(pi,n)
F	3	0	0
T	0	2	0

$$\checkmark \text{gain(windy)} = 0.97$$



TID	Items
T100	1, 3, 4
T200	2, 3, 5
T300	1, 2, 3, 5
T400	2, 5

minsup = 0.5

minconf = 50%

STEP 1

1(a) First Scan (item:count)

$$C_1 \rightarrow \{\{1\} : \frac{2}{4}, \{2\} : \frac{3}{4}, \{3\} : \frac{3}{4}, \{4\} : \frac{1}{4}, \{5\} : \frac{3}{4}\} \rightarrow \frac{1}{4} < 0.5 \text{ so discard}$$

$$F_1 \rightarrow \{1\} : \underline{2}, \{2\} : \underline{3}, \{3\} : \underline{3}, \{5\} : \underline{3}$$

(b) joining $C_2 \rightarrow$ make pairs of F_1

$$C_2 = (1,2), (1,3), (1,5), (2,3), (2,5), (3,5)$$

(c) pruning C_2

$$C_2 = \text{all items in } \underline{C_2} \text{ belong to } F_1$$

2(a) Second Scan (item:count)

$$C_2 = \{\{1,2\} : \frac{1}{4}, \{1,3\} : \frac{2}{4}, \{1,5\} : \frac{1}{4}, \{2,3\} : \frac{2}{4}, \{2,5\} : \frac{3}{4}, \{3,5\} : \frac{2}{4}\}$$

$$F_2 = \{1,3\} : \underline{2}, \{2,3\} : \underline{2}, \{2,5\} : \underline{3}, \{3,5\} : \underline{2}$$

(b) joining $C_3 \rightarrow$ make triples of F_2

$$C_3 = \begin{matrix} [1,3] \rightarrow [2,3] \\ [2,5] \\ [3,5] \end{matrix} \quad \begin{matrix} [2,3] \rightarrow [2,5] \\ [3,5] \end{matrix} \quad \begin{matrix} [2,5] \rightarrow [3,5] \end{matrix}$$

$$(c) \text{ pruning } C_3 = [2,3,5]$$

$$F_2 = [2,3], [2,5], [3,5] \text{ all belong to } C_3 = [2,3,5]$$

3(a) Third Scan (item:count)

$$C_3 = \{2,3,5\} : \frac{2}{4}$$

$$F_3 = \{2,3,5\} : 2$$

(b) joining $C_4 \rightarrow$ only 1 set of F_3 exist

$$C_4 = [2,3,5]$$

(c) pruning C_4

$$C_4 = [2,3,5]$$

STEP 2

$$\text{minconf} = 50\%$$

non empty subsets = $\{2,3\}, \{2,5\}, \{3,5\}, \{2\}, \{3\}, \{5\}$

$$60\% \quad 75\% \quad 50\% \quad 75\% \quad 75\% \quad 75\%$$

association rules :

$$\{2,3\} \rightarrow 5 \Rightarrow 2/2 = 100\% \quad \{2,3,5\}/\{2,3\}$$

$$\{2,5\} \rightarrow 3 \Rightarrow 2/3 = 67\% \quad \{2,3,5\}/\{2,5\}$$

$$\{3,5\} \rightarrow 2 \Rightarrow 2/2 = 100\% \quad \{2,3,5\}/\{3,5\}$$

$$\{2\} \rightarrow \{3,5\} \Rightarrow 2/3 = 67\% \quad \{2,3,5\}/\{2\}$$

$$\{3\} \rightarrow \{2,5\} \Rightarrow 2/3 = 67\% \quad \{2,3,5\}/\{3\}$$

$$\{5\} \rightarrow \{2,3\} \Rightarrow 2/3 = 67\% \quad \{2,3,5\}/\{5\}$$

AdaBoost

day / date:

CGPA	Interactivity	Practical Knowledge	Communication Skills	Job Profile
≥ 9	Yes	Good	Good	Yes
< 9	No	Good	Moderate	Yes
≥ 9	No	Average	Moderate	No
< 9	No	Average	Good	No
≥ 9	Yes	Good	Moderate	Yes
≥ 9	Yes	Good	Moderate	Yes

1. CGPA	C GPA	Predicted values	Actual values	Job Weight
	≥ 9	Yes	Yes	$\frac{1}{6} \cdot 0.167$
	< 9	No	Yes	$\frac{1}{6}$
	≥ 9	Yes	No	$\frac{1}{6}$
	< 9	No	No	$\frac{1}{6}$
	≥ 9	Yes	Yes	$\frac{1}{6}$
	≥ 9	Yes	Yes	$\frac{1}{6}$

2. calculate error and updated weights.

$$\text{correct} = 4 \times \frac{1}{6} \quad \left. \begin{array}{l} \\ \end{array} \right\} \quad \alpha = \frac{1}{2} \ln \left(\frac{1 - \frac{2}{6}}{\frac{2}{6}} \right) = 0.347.$$

$$\text{incorrect} = 2 \times \frac{1}{6}.$$

3. normalize all the weights.

$$\Rightarrow \frac{4 \times 0.167 \times e^{-0.347}}{0.472 + 0.945} + \frac{0.167 \times 2 \times e^{0.347}}{0.472 + 0.945} = \frac{\text{weight of correct} \times e^{-\infty}}{\text{weight of correct} \times e^{\infty} + \text{weight of incorrect} \times e^{\infty}}$$

$$\Rightarrow 0.472 + 0.945 = 0.9428.$$

4. calculate updated weight for correct and incorrect.

$$\text{new-weight-correct} = \frac{\frac{1}{6} \times e^{-0.34}}{0.9428} = 0.125$$

$$\text{new-weight-incorrect} = \frac{\frac{1}{6} \times e^{0.34}}{0.9428} = 0.248$$



KAGHAZ
www.kaghazpk

day / date:

2. Intervativeness

Intervativeness	Predicted values	Actual values	Weight
Yes	Yes	Yes	0.125
No	No	Yes	0.248
No	No	No	0.248
No	No	No	0.125
Yes	Yes	Yes	0.125
Yes	Yes	Yes	0.125

3 calculate error and updated weights

$$\text{correct} = 5 \quad \left. \begin{array}{l} \\ \end{array} \right\} \quad \alpha = \frac{1}{2} \ln \left(\frac{1 - 0.248}{0.248} \right) = 0.555$$
$$\text{incorrect} = 1$$

4. normalize all weights

$$= 0.125 \times 2^0 \times e^{-0.555} + 0.248 \times 1 \times e^{-0.555} + 0.248 \times 1 \times e^{0.555}$$
$$= 0.866$$

5. updated weights

$$\text{new-weight-correct} = \frac{0.125 \times e^{-0.555}}{0.866} = 0.0829$$

$$\text{new-weight-correct} = \frac{0.248 \times e^{-0.555}}{0.866} = 0.164$$

$$\text{new-weight-incorrect} = \frac{0.248 \times e^{0.555}}{0.866} = 0.499$$

3. Practical Knowledge

Practical knowledge	Predicted	Actual	Weights
Good	Yes	Yes	0.0829
Good	Yes	Yes	0.499
Average	No	No	0.164
Average	No	No	0.0829
Good	Yes	Yes	0.0829
Good	Yes	Yes	0.0829

4. calculate error and update weights.

correct : 6 $\alpha = \frac{1}{2} \ln \left(\frac{1-0}{0} \right)$

incorrect : 0

• weights will not be updated.

4. Communication Skill.

Communication Skill	Predicted	Actual	Weights
Good	Yes	Yes	0.0829 →
Moderate	No	Yes	0.499 ×
Moderate	No	No	0.164
Good	Yes	No	0.0829 ×
Moderate	No	Yes	0.0829 ×
Moderate	No	Yes	0.0829 +

5 calculate error and update weights

correct = 2 } $\alpha = \frac{1}{2} \ln \left(\frac{1-0.7477}{0.7477} \right) = -0.54$

incorrect = 4

? error = $0.499 \times 1 + 0.0829 \times 3 = 0.7477$



day / date:

b. normalize weights

$$= 0.0829 \times e^{-(-0.54)} + 0.164 \times e^{-(-0.54)} + 0.0829 \times 3 \times e^{-(0.54)}$$
$$+ 0.499 \times \bar{e}^{0.54}$$
$$= \frac{1.7067}{1.7067 + 0.7627 + 0.859} = 0.859$$

$$\frac{0.0829 \times e^{0.54}}{1.7067 + 0.7627 + 0.859} = \frac{0.0829}{1.7067 + 0.7627 + 0.859} \times 0.859 = 0.1663$$

$$\frac{0.164 \times e^{0.54}}{1.7067 + 0.7627 + 0.859} = \frac{0.1648}{1.7067 + 0.7627 + 0.859} \times 0.859 = 0.3331$$

$$\frac{0.0829 \times \bar{e}^{0.54}}{1.7067 + 0.7627 + 0.859} = \frac{0.0633}{1.7067 + 0.7627 + 0.859} = -0.0555$$

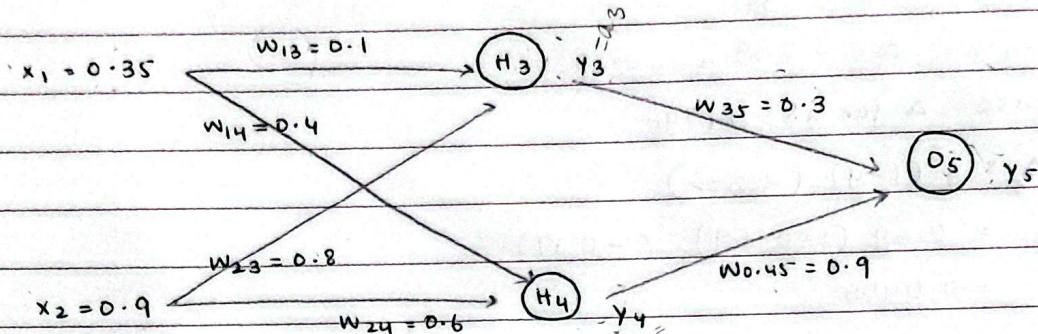
$$\frac{0.499 \times \bar{e}^{0.54}}{1.7067 + 0.7627 + 0.859} = \frac{0.381}{1.7067 + 0.7627 + 0.859} = 0.3337$$

→ conclusion	$\alpha' = 0.347$	$\alpha = 0.555$	$\alpha = -0.54$	weight	FP
Y	Y	Y	Y	0.362	Y
N	N	N	N	0	N
Y	N	N	N	0.347	Y
N	N	Y	Y	-0.45485	N
Y	Y	N	N	0.896	Y
Y	Y	N	N	0.896	Y



KAGHAZ
www.kaghaz.pk

BACK PROPAGATION


 $H_3 \rightarrow H_4 \rightarrow D_5$
 \downarrow
 a_3

$$\bullet \quad a_3 = (0.35)(0.1) + (0.8)(0.9)$$

sigmoid function = $\frac{1}{1+e^{-a_3}}$

$$SF = \frac{1}{1+e^{-0.8755}} = 0.68$$

$$\bullet \quad a_4 = (0.35)(0.4) + (0.9)(0.6) = 0.68$$

$$SF = \frac{1}{1+e^{-0.68}} = 0.6637$$

$$\bullet \quad y_{5as} = (0.68)(0.3) + (0.6637)(0.9) = 0.801$$

$$SF = \frac{1}{1+e^{-0.801}} = 0.69$$

day / date:

2. calculate error.

$$\text{error} = \text{target} - y_5$$

$$\text{error} = 0.5 - 0.69 = -0.19$$

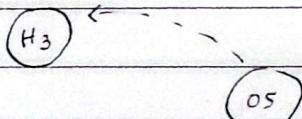
3. calculate Δ for y_3, y_4, y_5

$$\Delta_5 = y(1-y)(y\text{error})$$

$$= 0.69(1-0.69)(-0.19)$$

$$= -0.0406$$

$$\Delta_3 =$$



$$\Delta_3 = y(1-y)(w_{35} \times \Delta_5)$$

$$= 0.68(1-0.68)(0.35 \times -0.0406)$$

$$= -0.00265$$

$$\Delta_4 = y(1-y)(w_{45} \times \Delta_5)$$

$$= 0.6637(1-0.6637)(0.9 \times -0.0406)$$

$$= -0.0082$$

4. new weight for each line, 45, 35, 13, 14, 23, 24

$$\text{new weight} = \eta \times \Delta \times y$$

$$\eta = 1$$

$$w_{45} = 1 \times -0.0406 \times 0.68 = -0.0269$$

$$\text{new weight + original} = 0.9 + (-0.0269) = 0.8731$$

$$w_{35} = 1 \times -0.0406 \times 0.68 = -0.0276$$

$$\text{new + original} = 0.3 + (-0.0276) = 0.2724$$

$$w_{13} = 1 \times -0.00265 \times 0.35 = -0.00287$$

$$\text{new + original} = 0.1 + 0.00287 = 0.0991$$

$$w_{14} = 1 \times -0.0082 \times 0.35 = 0.3971 - 0.0287$$

$$\text{new + original} = 0.14 + (-0.0287) = 0.3971$$

$$w_{23} = 1 \times -0.00265 \times 0.9 = -0.002385$$

$$\text{new + original} = -0.002385 + 0.8 = 0.7976$$

$$w_{24}$$

$$\text{new + original} = 0.5926$$



KAGHAZ
www.kaghaz.pk