

01/01/22 Data warehouse and datamining

Data warehouse: collecting or gathering data from different sources and organized in centralized place is called data warehouse.

Ques → online analytical process

01/01/22

Mod-2 Data Preprocessing

→ The problem of making the data more suitable for data mining:

→ The tasks employed in this problem are informed by the process of data understanding.

(Washing the vegetables → removing unwanted things & considering only essential things)

meaningful

Extracting the insights from the data which drives decision-driven making

↓ Datamining

Increase the profit & minimize the loss

To extract the insights from data → Preprocess in first.

What is data?

- Collection of data objects and their attributes

- (a) Collection of raw facts information.

Each instance → structured objects

column → Attributes/features

- An attribute is a property or characteristic of an object.

- ex: eye color of a person, temperature etc.

- Attr is also known as parameters

Variable, field, characteristic, feature

- A collection of attributes describe an object in a particular way

- Object is also known as

record / point / case / example

entity / instance

missing value

- Absece of information in record → harmful consequence

to the validity of the subsequent analysis

- No missing values → Complete in records

- Missing values in incomplete record cases

Possibility of data correctness

✗

→ A value can't be -ve

Properties of Attribute Values

Attribute values

↓
are numbers or symbols
assigned to an attribute.

Distinction b/w attributes and attribute values:

Same attribute can be mapped to diff. attr. values.

E.g. height can be measured in both feet or meter.

Types of Attributes

$$y =, \neq$$

○ Nominal (fixed set of values)

ID numbers, eye color, [distinctness]
zip codes, Blender (MF)

○ Ordinal: (in ordering sense) [distinctness + order]

rankings (tarte & potato).

depends on a scale from 1-10,

grades, height in tall, medium,
shortly. ↗ ordering

○ Interval:

[distinctness, order & addition]

Calendar dates, temperature in

celcius & Fahrenheit

○ Ratio [all 4 properties]

temp in kelvin, length, time,

amount of money, weight of a

country, number of

1) Distinctness = \neq

2) Order $<$ $>$

3) Addition

4) Multiplication

X. eye color = black

Y. eye color = brown

X. eye color \neq Y. eye color

X \neq Y

X > Y \neq Y > X

X + Y \neq Y + X

X * Y \neq Y * X

ask in exam: specify type of attribute & properties

attribute given dataset

what for transforming over models

examples for ordinal, interval, ratio

examples for nominal, ordinal

examples for ratio

examples for interval

examples for ordinal, ratio

examples for nominal, ordinal

examples for ratio

examples for interval

examples for ordinal, ratio

examples for nominal, ordinal

examples for ratio

examples for interval

examples for ordinal, ratio

10/10/2022

Direct Attribute:

- countable set of values.

zip codes, country

- integer values.

- binary attr → apl conv.

Continuous Attribute:

- deal no.

- temp, height, weight

- It is Measured

continuous data → continuous domain

Record

data matrix

- Tables

- document data

each doc → term

each term → component

each comp → occurrence

each comp → no. times

wrt term occur.

frequency of words

No. of occurrences of words in documents.

Transaction data

Graph data

Node → no. of users using social media

Vertices → friend links.

Chemical data. C₆H₆



Ordered data

Sequence of transaction items every
(AB) (D) (CF)

Stream of data

Genomic sequence data.

A G C T

Adenine Guanine Cytosine Thymine

Spatio-temporal data

Avg. monthly temp.

Data quality problem

Noise or Outliers

Missing values

Duplicate data

Outliers

data objects with characteristics

considerably diff from majority
other data objects in the
data set.

Missing values (Absence of data)

Replaced by possible values

Eliminate data objects

Estimate missing value

Ignore it

Reason

- Attr. may not be

applicable to all cases

(salary for children X)

- Info not collected.

Duplicate data

Redundant data record.

Sample dataset: $\{4, 8, 9, 15, 21, 21, 24, 25, 26, 28, 34\}$

Instance:

feature:

Value 11:

redundant value → remove redundant

Data Preprocessing Tasks

- Data cleaning (removing unwanted / noisy data)
- Data Transformation (changing data representation)
- Data reduction (Reducing dimension)
- Data Discretization

✓ Numerical to nominal datatype conversion

Dividing the data

discrete bins, intervals

→ for better model performance

→ can be used for classification, regression,

clustering

Method of discretization: binning

→ equal width binning, equal frequency binning

Data Smoothing

$\{4, 8, 9, 15, 21, 21, 24, 25, 26, 28, 34\}$

equal frequency (equi-depth) binning:

Bin 1: $\{4, 8, 9, 15\}$

2: $\{21, 21, 24, 25\}$

3: $\{26, 28, 29, 34\}$

Smoothing by bin mean: (mean value)

Bin 1: $\{9, 9, 9, 9\}$

2: $\{23, 23, 23, 23\}$

3: $\{29, 29, 29, 29\}$

Smoothing by bin boundary

except min, max values; all are replaced

Bin 1: $\{4, 4, 4, 15\}$

2: $\{21, 21, 25, 25\}$

3: $\{26, 26, 26, 34\}$

110122-1ab

{ @RELATION <name>,
@ATTRIBUTE <attr-name> CATE,
@ATTRIBUTE class {< >},
@ATTRIBUTE class → Nominal

@ATTRIBUTE class {< >},
@ATTRIBUTE class → Nominal

13/10/22 Mean with outlier m_1
 Mean without outlier m_2 $m_1 - m_2 \neq 0$
 Outlier \rightarrow datapoint not close to other data.

Technique: IQR / boxplot used to detect outliers values. \rightarrow equivalent zero.

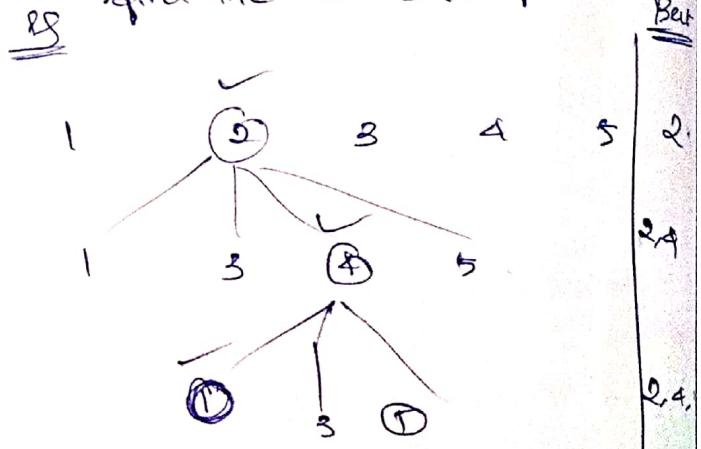
Feature selection

Sequential Forward Selection

Find the feature f_1 that gives the best performance.

Find the feature f_2 such that (f_1, f_2) gives the best perform repeat for as many feature as desire.

find the best 3/5 feature.

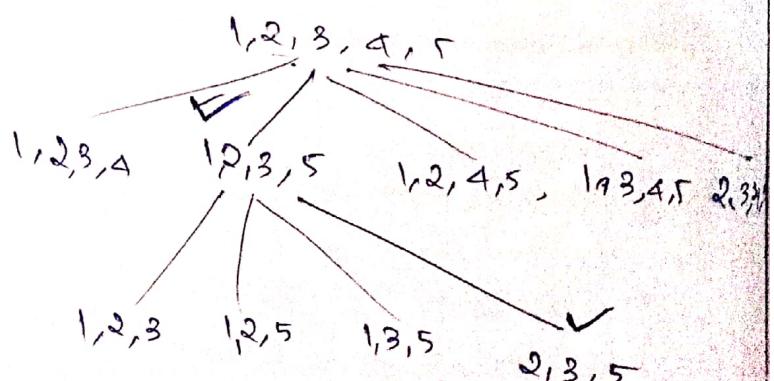


Sequential Backward Elimination

Buses	Miles
6 3	28
9	x

$$\frac{9}{x} = \frac{63}{28}$$

$$\boxed{x=9}$$



Data reduction

Age from 47 individuals

17, 18, 19, 20, 21, 22, 23, 24, 25, 26, 27, 28, 29, 30, 31, 32, 33, 34, 35, 36, 37, 38, 39, 40, 41, 42, 43, 44, 45, 46, 47

Data discretization in data mining

→ converts a large no. of data

values into smaller ones, so that

data evaluation & data management

becomes easy.

Age → 10, 11, 13, 14, 17, 19, 30,

32, 33, 35, 36, 37, 38, 39, 40, 41, 42, 43, 44, 45, 46, 47

Young, Mature, Old.

Binization

Mapping a categorical attribute
to a set of attribute that are
binary.

Ordinal attribute:

{Small, Medium, Large} →

00, 01, 10, 11



(old take any 3/4 to represent)

Qn. How

Ans 3

$$\log_2(3) = \text{ceil } \log_2(3) = \overline{\log_2(3)} = 3$$

Issues

x_2 & x_3 are now correlated

bz 'good' is encoded w/ both

Actual	0	0	0
poor	1	0	0
OK	2	0	1
good	3	0	1
Great	4	1	0

5 categorical values → 3 binary val.

Data reduction:-

Age of 47 individuals

$$\text{Range} = 69 - 20 = 49$$

$$N = 5 \text{ (No. of intervals)}$$

$$\text{Width of class} = 49/5 = 9.8 = 10$$

$$\text{class interval} = 20 - 29, 30 - 39,$$

$$40 - 49, 50 - 59, 60 - 69$$

$$70 - 79, 80 - 89, 90 - 99$$

Continuous: 20, 21, 22, .. 69

Discrete: 20, 21, 22, .. 69

Interval: 20-29, 30-39

↓

Ordinal: two-tier, three

↓

Nominal: Young, Middle,

not dots reduce

data
reduction

ordinal → discrete

continuous → discrete

nominal → discrete

interval → discrete

continuous → continuous

ordinal → continuous

nominal → continuous

interval → continuous

continuous → continuous

ordinal → continuous

nominal → continuous

interval → continuous

continuous → continuous

ordinal → continuous

nominal → continuous

interval → continuous

continuous → continuous

ordinal → continuous

nominal → continuous

interval → continuous

continuous → continuous

ordinal → continuous

nominal → continuous

interval → continuous

continuous → continuous

ordinal → continuous

nominal → continuous

interval → continuous

continuous → continuous

ordinal → continuous

nominal → continuous

interval → continuous

continuous → continuous

ordinal → continuous

nominal → continuous

interval → continuous

continuous → continuous

ordinal → continuous

nominal → continuous

interval → continuous

continuous → continuous

ordinal → continuous

nominal → continuous

interval → continuous

continuous → continuous

ordinal → continuous

nominal → continuous

interval → continuous

continuous → continuous

ordinal → continuous

nominal → continuous

interval → continuous

continuous → continuous

ordinal → continuous

nominal → continuous

interval → continuous

continuous → continuous

ordinal → continuous

nominal → continuous

interval → continuous

continuous → continuous

ordinal → continuous

nominal → continuous

interval → continuous

continuous → continuous

ordinal → continuous

nominal → continuous

interval → continuous

continuous → continuous

ordinal → continuous

nominal → continuous

interval → continuous

continuous → continuous

ordinal → continuous

nominal → continuous

interval → continuous

continuous → continuous

ordinal → continuous

nominal → continuous

interval → continuous

continuous → continuous

ordinal → continuous

nominal → continuous

interval → continuous

continuous → continuous

ordinal → continuous

nominal → continuous

interval → continuous

continuous → continuous

ordinal → continuous

nominal → continuous

interval → continuous

continuous → continuous

ordinal → continuous

nominal → continuous

interval → continuous

continuous → continuous

ordinal → continuous

nominal → continuous

interval → continuous

continuous → continuous

ordinal → continuous

nominal → continuous

interval → continuous

continuous → continuous

ordinal → continuous

nominal → continuous

interval → continuous

continuous → continuous

ordinal → continuous

nominal → continuous

interval → continuous

continuous → continuous

ordinal → continuous

nominal → continuous

interval → continuous

continuous → continuous

ordinal → continuous

nominal → continuous

interval → continuous

continuous → continuous

ordinal → continuous

nominal → continuous

interval → continuous

continuous → continuous

ordinal → continuous

nominal → continuous

interval → continuous

continuous → continuous

ordinal → continuous

nominal → continuous

interval → continuous

continuous → continuous

ordinal → continuous

nominal → continuous

interval → continuous

continuous → continuous

ordinal → continuous

nominal → continuous

interval → continuous

continuous → continuous

ordinal → continuous

nominal → continuous

interval → continuous

continuous → continuous

ordinal → continuous

nominal → continuous

interval → continuous

continuous → continuous

ordinal → continuous

nominal → continuous

interval → continuous

continuous → continuous

ordinal → continuous

nominal → continuous

interval → continuous

continuous → continuous

ordinal → continuous

nominal → continuous

interval → continuous

continuous → continuous

ordinal → continuous

nominal → continuous

interval → continuous

continuous → continuous

ordinal → continuous

nominal → continuous

interval → continuous

continuous → continuous

ordinal → continuous

nominal → continuous

interval → continuous

continuous → continuous

ordinal → continuous

nominal → continuous

interval → continuous

continuous → continuous

ordinal → continuous

nominal → continuous

interval → continuous

continuous → continuous

ordinal → continuous

nominal → continuous

interval → continuous

continuous → continuous

ordinal → continuous

nominal → continuous

interval → continuous

continuous → continuous

ordinal → continuous

nominal → continuous

interval → continuous

continuous → continuous

ordinal → continuous

nominal → continuous

interval → continuous

continuous → continuous

ordinal → continuous

nominal → continuous

interval → continuous

continuous → continuous

ordinal → continuous

nominal → continuous

interval → continuous

continuous → continuous

ordinal → continuous

nominal → continuous

interval → continuous

continuous → continuous

ordinal → continuous

nominal → continuous

interval → continuous

continuous → continuous

ordinal → continuous

nominal → continuous

interval → continuous

continuous → continuous

ordinal → continuous

nominal → continuous

interval → continuous

continuous → continuous

ordinal → continuous

nominal → continuous

interval → continuous

continuous → continuous

ordinal → continuous

nominal → continuous

Data Transformation: Normalization:

- Min-max
- Z-score
- by decimal scaling

Min-max (new_min , new_max)

$$v' = \frac{v - \text{min}}{\text{max} - \text{min}} (\text{new_max} - \text{new_min}) + \text{new_min}$$

let income range

from \$12,000 to \$98,000 normalized to

range of decimal scaling from (0.0 - 1.0).

\$73,600 is mapped to

$$\frac{73,600 - 12,000}{98,000 - 12,000} (1 - 0) + 0 = 0.716$$

68,000

$$\frac{68,000 - 12,000}{98,000 - 12,000} / 10 = 0.6511$$

$$= \frac{56,000}{86,000} \times 10 = 0.6511$$

Z-score (μ : mean, σ : std)

$$v' = \frac{v - \mu}{\sigma}$$

$$\Sigma \quad \mu = 54,000$$

$$\sigma = 16000$$

$$\text{then } \frac{73,600 - 54,000}{16,000} = 1.225$$

Normalization by decimal scaling:

$$v' = \frac{v}{10^j} \quad j \rightarrow \text{smallest integer such that } \text{Max}(v') < 1$$

22/01/21

Correlation

	Harm(H)	Mast (m)
data	9	39
	15	56
	25	93
	14	61
	10	50
	18	75
	6	32
	16	85
	5	42
	19	70
	16	66
	26	80
total	161	749
Avg:	13.92	62.42

$$\text{Var}(x) = \frac{\sum_{i=1}^n (x_i - \bar{x})(x_i - \bar{x})}{(n-1)}$$

$$\text{cov}(x, y) = \frac{\sum_{i=1}^n (x_i - \bar{x})(y_i - \bar{y})}{n-1}$$

positive $\rightarrow x \uparrow y \uparrow, x \downarrow y \downarrow$ negative $\rightarrow x \uparrow y \downarrow, x \downarrow y \uparrow$ sp. $x \uparrow y \uparrow, x \downarrow y \downarrow$

H	M	(H _i - H̄)	(M _i - M̄)	(H _i - H̄)(M _i - M̄)
9	39	-4.92	-23.42	115.28
15	56	1.08	-6.42	-6.93
25	93	11.08	30.58	338.83
14	61	0.08	-12.42	-0.11
10	50	-3.92	12.58	48.69
18	75	4.08	-30.42	51.33
6	32	-13.92	22.58	423.45
16	85	2.08	-20.42	46.97
5	42	-8.92	7.58	182.15
19	70	5.08	3.58	38.51
16	66	2.08	1.58	7.45
26	80	6.08		166.89
				1149.89
tot.				104.54

Avg:

$$\text{cov}_{xy} = \frac{\sum xy - \frac{\sum x \sum y}{n}}{n-1}$$

$$\gamma = \frac{\text{cov}_{xy}}{s_x s_y}$$

$$\gamma = \frac{-1}{\sqrt{n-1}}$$

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

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

Forward feature selection

1/1 Implementation

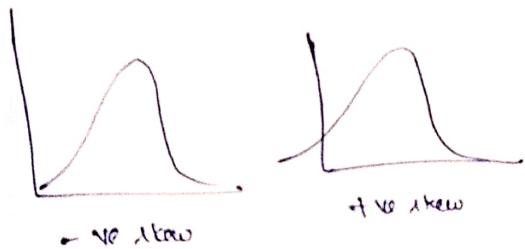
Variable transformation

In data modeling, transformation refers to the replacement of a variable by a function.

for instance, replacing a variable by

square root or $\log x$ is a transformation

Transformation is a process that changes the distribution or relationship of a variable with others.



vs skew

↓
transformation



(10,12) (2,2)

$\text{Sim}(A, B)$

$$\frac{10 \times 2 + 12 \times 2}{\sqrt{100+104} \times \sqrt{4+4}} = \frac{20+24}{\sqrt{204} \times \sqrt{8}}$$

$$= \frac{44}{44.18} \approx 1$$

Similarity vs Dissimilarity

Similarity

Numerical Measure

of how alike two data obj are.

- is higher when obj are more alike
- often falls in range [0, 1]

$$\text{Sim}(x_1, x_2) = \frac{0.9}{0.3}$$

based on similarity consider any attr. if both are similar.

Dissimilarity

- Numerical measure of how diff. two are & data obj!
- lower than objects are more alike
- Minimum dissimilarity often 0.
- Upper limit varies.

Proximity refers to a similarity or dissimilarity

$$\text{dissimilarity} = 1 - \text{similarity}$$

$$\text{Similarity } (A|B) = \frac{A \cdot B}{\|A\| \|B\|} = \frac{\sum_{i=1}^n A_i B_i}{\sqrt{\sum_{i=1}^n A_i^2} \times \sqrt{\sum_{i=1}^n B_i^2}}$$

$$\text{Similarity } (x,y) = \cos \theta = \frac{x \cdot y}{\|x\| \|y\|}$$

Atr-Hpo

disimilarity

Nominal

$$d = \begin{cases} 0 & \text{if } p=q \\ 1 & \text{if } p \neq q \end{cases}$$

Ordinal

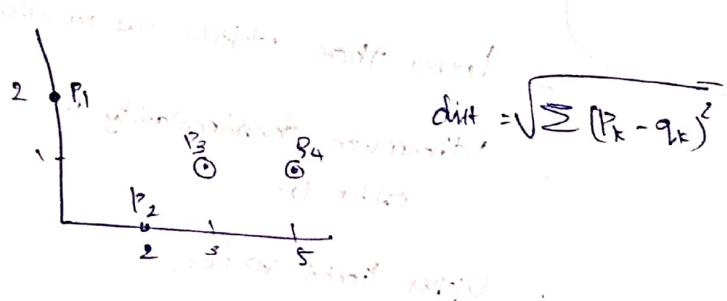
$$d = \frac{|p-q|}{n-1}$$

(values mapped to integers 0 to n-1,
n-1 no p value)

$$d = |p-q|$$

Interval & Ratio

Euclidean distance



Point

	x	y
P ₁	0	2
P ₂	2	0
P ₃	3	1
P ₄	5	1

Distances

	P ₁	P ₂	P ₃	P ₄
P ₁	0	2.828	3.162	5.099
P ₂	2.828	0	1.414	3.162
P ₃	3.162	1.414	0	2
P ₄	5.099	3.162	2	0

Distance matrix

Similarty

$$S = \begin{cases} 1 & \text{if } p=q \\ 0 & \text{if } p \neq q \end{cases}$$

$$S = \frac{1}{n-1} \sum_{i=1}^n |p_i - q_i|$$

or get standard p, normalized 2nd

$$S = -d, S = \frac{1}{n-1} \sum_{i=1}^n |p_i - q_i|$$

$$S = 1 - \frac{d - \min d}{\max d - \min d}$$

min d = min distance

max d = max distance

min d = min distance

max d = max distance

min d = min distance

max d = max distance

min d = min distance

max d = max distance

min d = min distance

max d = max distance

min d = min distance

max d = max distance

min d = min distance

max d = max distance

min d = min distance

max d = max distance

min d = min distance

max d = max distance

min d = min distance

max d = max distance

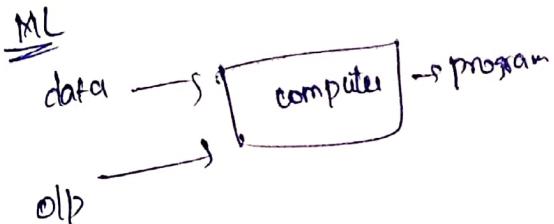
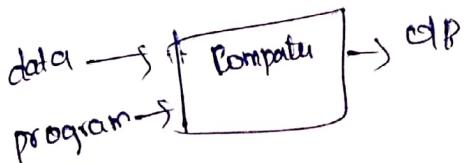
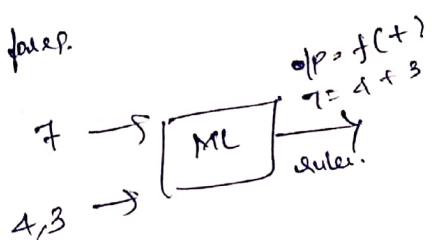
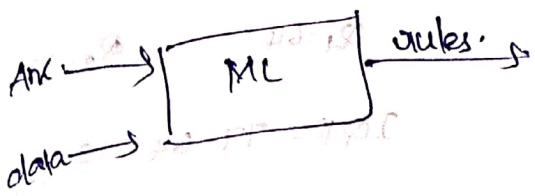
Not Regular programming

Field of study that gives computer

ability to learn without being explicitly programmed.
— Arthur Samuel (1959)



Machine Learning



ML - Machine Learning :-

Learning is any process by which a

system improves performance from experience

— Herbert A. Simon

ML is the study of algorithms that

improve their performance at some tasks.

Types of learning :-

Supervised (Inductive) learning

training data + desired o/p (labels)

Unsupervised learning

training data (without desired o/p)

Semi-supervised learning

training data + a few desired o/p

Reinforcement learning

Rewards from sequence of actions.

(Q, R)

23/10/2023

Date

No. of animal distract in each kid's box

(B) $4, 4, 6, 7, 11, 12, 14, 15, 6$ $IQR = ?$

Sorting

~~4, 4, 6, 7, 11, 12, 14, 15,~~

$$\text{Median}, Q_2 = \frac{N}{2} = 9$$

$\approx 9 \frac{1}{2} = 15^{\text{th}}$ element.

$$Q_2 = 10$$

~~4, 4, 6, 7~~

$$Q_1 \Rightarrow \text{Median of first set} = \frac{4+6}{2} = \frac{10}{2}$$

$$Q_1 = 5$$

11, 12, 14, 15

$$Q_3 \Rightarrow \text{Median} = \frac{12+14}{2} = \frac{26}{2} = 13$$

$$Q_3 = 13$$

$$IQR = Q_3 - Q_1$$

$$= 13 - 5 = 8$$

$$IQR = 8$$

Upper fence \Rightarrow

$$Q_3 + 1.5 \times IQR$$

~~Q₃ + 1.5 × IQR~~

$$= 13 + 1.5 \times 8$$

$$= 13 + 12 = 25$$

Lower \Rightarrow

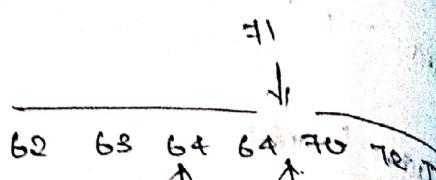
$$Q_1 - 1.5 \times IQR$$

$$5 - 1.5 \times 8$$

$$5 - 12 = -7$$

$$\text{lower} = -7$$

②



$$IQR = 77 - 64 = 13$$

Q₁ = 64

Q₃ = 70

37/01/2022

ML

inp x out y

bike label 'motorcycle' (Supervised)

Classification:-

Classification is a process of categorizing a given set into classes, it can be performed on both structured or unstructured data. The process starts with predicting the class of given data points. The classes are often referred to as target, label, categories.

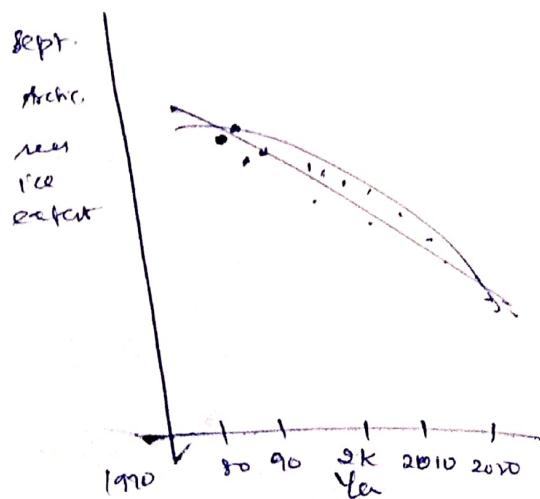
Ex. Mail spam/not spam.

human decision making - classes

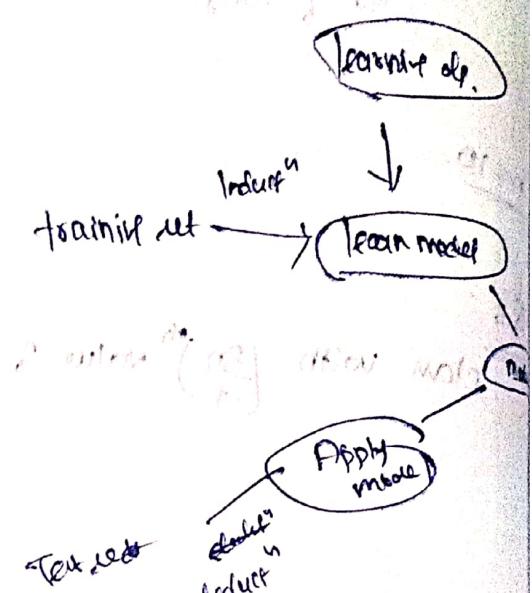
Supervised Learning: Regression

$f(x)$ is predict y given x

y is real valued = regression



$$y = mx + c$$



Evaluation 2 classification models

Count of test records → correctly predicted by the classifier "model".
 Confusion matrix

		Predicted class	
		Class = 1	Class = 0
Actual class	Class = 1	f_{11}	f_{10}
	Class = 0	f_{01}	f_{00}

		Actual	
		Positive	Negative
Predicted	Positive	TP	FP
	Negative	FN	TN

$$\text{Accuracy} = \frac{\# \text{ correct predictions}}{\text{total } \# \text{ of predictions}}$$

$$= \frac{f_{11} + f_{00}}{f_{11} + f_{10} + f_{01} + f_{00}}$$

$$\text{Error rate} = \frac{\# \text{ wrong predictions}}{\text{total } \# \text{ of predictions}}$$

$$= \frac{f_{10} + f_{01}}{f_{11} + f_{10} + f_{01} + f_{00}}$$

$$\text{Accuracy} = \frac{(TN + TP)}{TN + FP + FN + TP}$$

$$\text{Precision} = \frac{TP}{TP + FP}$$

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

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

Measure Formula

$$\text{Recall} = \frac{TP}{P}$$

$$\text{F measure} = \frac{2 \times \text{Precision} \times \text{Recall}}{(\text{Precision} + \text{Recall})}$$

Area

Recall (Sensitivity)

(Sensitivity)

True positive recall

$$\frac{TP}{TP + FN}$$

		Actual	
		TP	FP
Predicted	TP	TP	FP
	FN	FN	TN

Precision

$$\text{Precision} \rightarrow \frac{TP}{TP + FP}$$

True Negative Rate

(Specificity)

$$\frac{TN}{TN + FP}$$

		Actual	
		TP	FP
Predicted	TP	TP	FP
	FN	FN	TN

False positive rate

$$\frac{FP}{FP + TN}$$

Actual

		Actual	
		TP	FP
Predicted	TP	TP	FP
	FN	FN	TN

False Positive Rate

(1 - Specificity)

$$\frac{FP}{FP + TN}$$

- Non-interesting file input

- bad output

- bad output

01/02/22

Decision Tree classifier:-

→ A "classification" scheme which generates a tree and a set of rules from given data set.

→ the set of records available for developing classification

method is divided into 2 disjoint subsets - training set

- testing set.

—

inner node → Attribute

edge → test on father node attr.

leaf → class.

— Topdown strategy

— Based on training data.

—

DT - ID3 algorithm

— Entropy [large]

— Gain [high value]

DT - CART

— Gini [small].

Naive Bayes Model

$$P(C_k|B) = \frac{P(C_k) \cdot P(x|C_k)}{P(x)}$$

$$P(A|B) = \frac{P(B|A) \cdot P(A)}{P(B)}$$

$P(x|C_k)$ is the posterior probability
of class k, given predictor x.

$P(C_k)$ → prior prob. of class k.

prior → prior prob. of predictors.

$P(x|C_k)$ → current likelihood of
predictor given class.

Ex

outlook = sunny, Temp = cool, Humidity = high,
wind = starts.

$$P(\text{PlayTennis} = \text{Yes}) = 9/14 = 0.64$$

$$P(\text{Play} = \text{No}) = 5/14 = 0.36$$

Outlook Y N

Sunny 9/9 3/5

Overcast 4/9 0

Rain 3/9 2/5

Temp Y N

Hot 2/9 2/5

Mild 4/9 2/5

Windy 3/9 1/5

Humidity Y N
High 3/9 4/5

Normal 6/9 1/5

Windy Y N

Strong 3/9 3/5

Weak 6/9 2/5

$$V_{NB} = \operatorname{argmax}_{V_j} P_{M_j} \prod_i P(a_i|v_j)$$

$$= \operatorname{argmax}_{V_j} P(V_j)$$

$$V_j \in \{\text{yes}, \text{no}\}$$

$$P(\text{outlook} = \text{sunny}|V_j) P(\text{Temp} = \text{cool}|V_j)$$

$$P(\text{Humidity} = \text{high}|V_j) P(\text{Wind} = \text{strong}|V_j)$$

$$V_{NB}(\text{Yes}) = P(\text{Yes}) P(\text{sunny|Yes})$$

$$P(\text{cool|Yes}) P(\text{high|Yes}),$$

$$P(\text{strong|Yes})$$

$$= 0.0053$$

$$V_{NB}(\text{No}) = P(\text{No}) P(\text{sunny|No})$$

$$P(\text{cool|No}) P(\text{high|No})$$

$$P(\text{strong|No})$$

$$= 0.0206$$

$V_{NB}(\text{No})$ is higher.

$$V_{NB}(\text{Yes}) = \frac{V_{NB}(\text{Yes})(\text{Yes})}{V_{NB}(\text{Yes}) + V_{NB}(\text{No})}$$

$$= 0.005$$

$$V_{NB}(\text{No}) = \frac{V_{NB}(\text{No})}{V_{NB}(\text{Yes}) + V_{NB}(\text{No})}$$

$$= 0.795$$

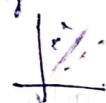
Shows $V_{NB}(\text{No})$ is high, so given no
rain is No.

Regression Analysis

Study of 2 variables in an attempt to find a relationship, or correlation.

Linear Regression $y = \beta_0 + \beta_1 x$

estimate β_0, β_1



x	y
-2	1
1	1
3	2

$$y = \alpha x + b$$

slope intercept

using the least square method

$$\alpha = \frac{1}{n} \sum_{i=1}^n x_i y_i - \frac{\sum x_i \sum y_i}{n}$$

using the least square method

$$n \sum x_i^2 - (\sum x_i)^2$$

$$b = \frac{1}{n} \left[\sum y_i - \alpha \sum x_i \right]$$

x	y	xy	x^2
-2	-1	2	4
1	1	1	1
3	2	6	9

$$\sum x = 2 \quad \sum y = 2 \quad \sum xy = 9 \quad \sum x^2 = 12$$

Calculate value of α

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

$$= \frac{3(9) - 2 \times 2}{3 \times 12 - 4^2} = \frac{24}{4} = 6$$

$$\alpha = \frac{23}{38}$$

$$b = \frac{1}{n} [\sum y - \alpha \sum x]$$

$$= \frac{1}{3} [2 - \frac{23}{38} \times 2]$$

$$= \frac{2}{3} \left[\frac{38 - 23}{38} \right]$$

$$= \frac{2}{3} \times \frac{15}{38}$$

$$b = \frac{5}{19}$$

$$y = \frac{23}{38} x + \frac{5}{19}$$