

DATA PREPROCESSING $\chi^2$  analysis

$$- \chi^2 = \sum (O - E)^2 \quad O - \text{Observed Value}$$

$\in$

$$E - \text{Expected Value.}$$

Calculation of E  $\rightarrow$  For ij in contingency table,  
 compute  $\frac{\text{rowsum}(i) \times \text{colsum}(j)}{\text{Total Elements}}$

- It is a test of independence for two attributes (categorical).
- It is a right tailed test.
- $H_0$ : The two attributes are independent/not correlated.
- dof:  $(r-1) \times (c-1)$  where (r, c) are rows & cols in contingency table  
 i.e. They refer to the levels in a category.
- $\chi^2 \leq \chi^2_{\text{critical}} \Rightarrow$  Fail to reject the Null Hypothesis (Independent)
- $\chi^2 > \chi^2_{\text{critical}} \Rightarrow$  Reject the null hypothesis. (Not Independent)

PEARSON'S CORRELATION CO-EFFICIENT

$$r_{xy} = \frac{\sum (x_i - \bar{x})(y_i - \bar{y})}{n \sigma_x \sigma_y} = \frac{\text{Cov}(x, y)}{\sqrt{\text{Var}(x) \cdot \text{Var}(y)}} = \frac{\sum x_i y_i - n \bar{x} \bar{y}}{\sqrt{n} \sqrt{\text{Var}(x) \cdot \text{Var}(y)}}$$

$r_{xy} > 0 \Rightarrow$  Positive correl

$< 0 \Rightarrow$  Negative correl

$r_{xy} \in [-1, 1]$

Sign of r  $\rightarrow$  Direction of relation

Magnitude of r  $\rightarrow$  Strength (linear) of relation

$r = 0 \Rightarrow$  No linear relation between x & y.

BINNING

$\rightarrow$  EQUAL DEPTH

Arrange data in asc. order; det. no. of bins

Dataset Size/Bin number  $\rightarrow$  Bin Capacity

Create Bins.

$\rightarrow$  EQUAL WIDTH

Arrange data in asc. order; det. min & max

Bin Width =  $(\text{max} - \text{min}) / \text{no. of bins}$

Formulate bins  $(x_i + w, x_i + 2w]$

where  $x_{i-1} \rightarrow$  Prev. bin close & w - bin width.

SMOOTHENING —

- By Mean → Replace all elems in a bin by mean of those elems
- By Median → " median = "
- By Boundaries → Replace elems closer to min of bin by min & closer to max of bin by max.

NORMALIZATION —

- Min Max Scaling →  $x' = \frac{x - x_{\min}}{x_{\max} - x_{\min}}$
- Decimal Scaling →  $x' = \frac{x}{10^j}$  → ( $j$  is the order of largest no. in dataset)

## DATA EXPLORATION

Attribute types

Qualitative → Nominal

→ Binary

Ordinal (Order matters)

Quantitative → Interval Scaled (Only S matters)

→ Ratio Scaled (Ht & Wt, contains abs. zero)

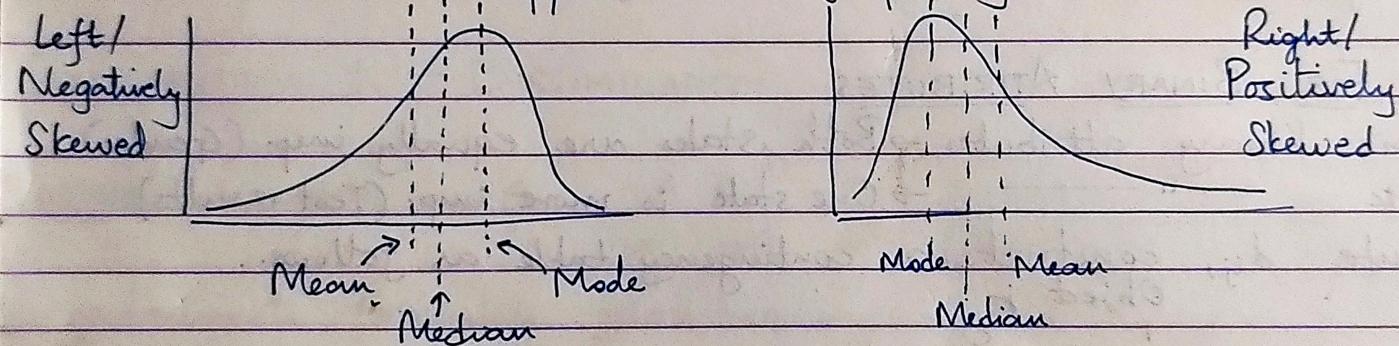
Discrete Attribute → Has finite or countably infinite set of values.

Continuous Attribute → Has Real numbers as attribute values.

Mean →  $\frac{\sum x_i}{N}$ ; Trimmed Mean  $\frac{\sum x_i}{N} \text{ } \forall x_i \notin \{\text{Outlier set}\}$

Median → Arrange data in ascending order; Middle element → Median.

Mode → Item which appears most frequently in the dataset.



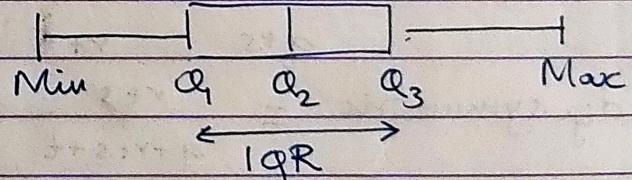
## FIVE NUMBER SUMMARY

Min  $Q_1$   $Q_2$  ( $Q_2$  Median)  $Q_3$  Max

$Q_1 \rightarrow$  Element @  $\frac{(n+1)}{4}$  position

$Q_2 \rightarrow$  Element @  $\frac{n+1}{2}$  position

## Box & Whisker's Plot



$Q_3 \rightarrow$  Element @  $\frac{3(n+1)}{4}$  position.

## VARIANCE

$$\text{Variance } (\sigma^2) = \frac{1}{n} \sum_{i=1}^n (x_i - \bar{x})^2$$

Population v/s Sample

$$\text{Variance } (s^2) = \frac{1}{n-1} \sum_{i=1}^n (x_i - \bar{x})^2$$

$$\text{Population stdv } \sigma = \sqrt{\sigma^2}$$

$$\text{Sample stdv } s = \sqrt{s^2} \rightarrow \text{Sqrt. of Sample Variance}$$

## Similarity / Dissimilarity (Proximity Measures)

Data matrix - Input data (with all attributes but the target)

Dissimilarity Matrix  $\rightarrow$  A  $\Delta$  matrix with  $d_{ij}$  = dissim. between record i & record j

$\rightarrow$  Technically it's a symmetric matrix with  $d_{ij} = 0 \forall i=j$

### Proximity For Nominal Attributes

$$d_{ij} = \frac{p-m}{p}$$

p - No. of nominal attributes

m - No. of matches

$$d_{ij} \in [0, 1]$$

\* If you're forced to weigh all attributes equally in a mixed attribute dataset, don't use this. Only compare match/no match.

### Proximity For Binary Attributes

Symmetric binary attributes  $\rightarrow$  Both states are equally imp. (Gender)

Asymmetric  $\rightarrow$  One state is more imp. (Test results)

To compute  $d_{ij}$ , construct a contingency table as follows.

Object j

		0	$q+r$	$q+r+s+t$
Object i	1	$q$	$r$	
	0	$s$	$t$	$s+t$

$$q+s \quad r+t$$

$$d_{ij, \text{symmetric}} = \frac{r+s}{q+r+s+t}$$

$$= \frac{\text{Items which don't match}}{\text{Total Items}}$$

$$d_{ij, \text{asymm.}} = \frac{r+s}{q+r+s}$$

$\hookrightarrow$  Don't care about (0,0) match for similarity.

### Proximity For Numeric Data

$$d_{ij} = \sqrt[n]{\sum (x_i - x_j)^n}$$

If  $n=2 \Rightarrow$  Euclidean

$n=\mathbb{R} - \{0, 1, 2\} \Rightarrow$  Minkowski

$n=\infty \Rightarrow$  Supremum  $\rightarrow$  Max value sum

$n=1 \Rightarrow$  L1-norm distance (Manhattan)

$p_1$	0	2	$(p_1, p_2)_1 =  2  +  2  = 2$
$p_2$	2	0	$(p_1, p_3)_2 = \sqrt{3^2 + 1^2} = \sqrt{10}$
$p_3$	3	1	$(p_1, p_4)_\infty = \max(5, 1) = 5$
$p_4$	5	1	

## PROXIMITY MEASURE OF ORDINAL DATA

→ Replace each ordinal by rank

→ Normalize the ranking  $r = \frac{r_{if} - 1}{M_f - 1}$

$r_{if}$  → Rating of record i for ordinal feature f

$M_f$  → Number of levels in feature f.

→ Compute the distance using Euclidean dist. formula

## MIXED DATATYPE

Find dissimilarity based on individual attribute type as we discussed so far & then weight them

$$d_{ij} = \frac{\sum w_f \cdot d_{ijf}}{\sum w_f} \rightarrow \begin{array}{l} w_f \rightarrow \text{Feature weight} \\ d_{ij} \rightarrow \text{Distance values.} \end{array}$$

SIMILARITY = 1 - DISSIMILARITY

COSINE SIMILARITY  $\rightarrow \cos \theta = \frac{A \cdot B}{\|A\| \|B\|} = \frac{\sum a_i b_i}{\sqrt{\sum a_i^2 \cdot \sum b_i^2}}$

Given documents formulate term frequency matrix

	$d_1$	$d_2$	$d_3$	$d_4$	...
$t_1$					}
$t_2$					
$t_3$					
$t_4$					
:					

Each column is now a vector for document

Each row is now a vector for terms

## CLASSIFICATION

### DECISION TREES

Algorithms: ID3, CART, C4.5 etc.

Principle → Split records based on an attribute that optimizes a criterion

- Splits -
- Binary Split → Straightforward for binary/categorical
  - Multi-way split → 2 bins for continuous
  - As many splits as levels in categorical
  - Discretize bins for continuous attributes

#### Measures of Node Impurity

$$\text{Gini Index} \rightarrow GI = 1 - \sum_j p_j^2 \quad \text{where } p_j = \frac{n_j}{n_{\text{dataset}} \text{ (node)}} \quad \& \quad r \text{ is number of target classes.}$$

$$= 1 - \sum_{j=1}^r \left( \frac{n_j}{n} \right)^2$$

$$\text{Gini Index}_{(\text{Node Split})} = \sum_i w_i \cdot \text{Gini}_i$$

Here  $w_i \rightarrow$  Proportion of  $i^{\text{th}}$  level in the node split

$\text{Gini}_i \rightarrow$  Corresponding Gini Index of that partition.

$$\Delta GI = \text{Gini}(\text{Node}) - \text{Gini}_{(\text{Node Split})}$$

For every attribute A, compute  $\Delta GI$  and select the one with highest value of  $\Delta GI \Rightarrow$  Lowest value of  $\text{Gini}_{(\text{Node Split})}$

$$\text{Entropy} (E) \quad E = - \sum p_j \log p_j \quad (\text{Everything else is same as above})$$

IG i.e. Information Gain = Entropy (Node) - Entropy (Node Split)

Select attribute with high IG for splitting at any given pt.

Decision boundaries are always // to the axes

If Split involves multiple attributes then boundaries can become oblique  
⇒ Oblique Decision Trees.

$$\text{Gain Ratio}_A = \frac{\text{Information Gain}(A)}{\text{Entropy}(\text{NodeSplitA})}$$

→ Attribute with maximum gain ratio is selected as splitting attribute.

## RULE-BASED CLASSIFIER

Rule (Condition)  $\rightarrow$  y  
Antecedent      Consequent

A rule  $r$  covers an instance  $x$  if attributes of  $x$  satisfy cond. in  $r$

Coverage  $\rightarrow$  Fraction of records which satisfy the rule antecedent

Accuracy  $\rightarrow$  \_\_\_\_\_ both antecedent & consequent

### Mutually Exclusive Rules

- If rules are independent of each other
- Every record is covered by at most 1 rule.

### Exhaustive Rules

- If it accounts for all possible combinations of attribute values.
- Each record is covered by at least 1 Rule.
- May lead to clash of  $\geq 1$  Rule
  - Ordered set
  - Majority Voting

### Rule Quality measures

$$\text{FOIL Gain} \rightarrow \text{Gain}(P_0, R_i) = t \left\{ \log \frac{P_1}{P_1 + u_i} - \log \frac{P_0}{P_0 + u_0} \right\}$$

## REGRESSION

$$Y = a + b \bar{x}$$

Dep |      | Indep  
 (Predicted) |      | (Predictor)  
 Y-intercept      Slope

LS method (Least Squares Sum)

$$b = \frac{\sum xy - n \bar{x} \bar{y}}{\sum x^2 - n \bar{x}^2}$$

$$a = \bar{y} - b \bar{x}$$

$$\text{Std Error } S_{xy} = \sqrt{\frac{\sum (y_i - y')^2}{n-2}}$$

## EVALUATION METRICS

~~Predicted →~~

<del>Actual ↓</del>	C1	C0	Recall = $\frac{TP}{TP+FN}$
C1	TP	FP	Sensitivity = Recall
C0	FN	TN	

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

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

$$F\text{-score} = \frac{2 \times \text{Precision} \times \text{Recall}}{\text{Precision} + \text{Recall}}$$

$$F_\beta = \frac{(1+\beta^2) \times \text{Precision} \times \text{recall}}{\beta^2 \times (\text{Precision}) + \text{recall}}$$