

Classification	
<ul style="list-style-type: none"> <li>- Split data accord. to some criterion, classify new (o) accord. to same criterion</li> <li>- P( new (o) belong to each group ) calc.</li> <li>- New (o) assigned to group with highest prob.</li> </ul>	
Discriminant Analysis (Supervised learning)	Cluster Analysis (Unsupervised learning)
Groups specified (can identify from available data)	Groups not specified, need to find from data itself Find hidden assoc. structures within unlabeled data
Logistic regression	Clustering
Training data: data used to build model Misclassification: predict wrongly (Goal: min. this error)	Survival analysis: Failures of Samsung phones
Quote prob. rather than binary 0-1 decision Labels of response variable known	Group similar objects but finds degree and nature of similarity by itself
	<ul style="list-style-type: none"> <li>- Preliminary classifier: rough labelling of (o) before thorough model-based analysis attempted</li> <li>- Determine groups for discriminant analysis</li> </ul>
	K-means clustering <ol style="list-style-type: none"> <li>1. Choose value of k, guess for k cluster centriods</li> <li>2. Compute distance of each pt from k centroids (Assign each pt. to cluster with closest centroid)</li> <li>3. Recompute centroids (cluster mean) for each updated cluster</li> <li>4. Con't. step 2 &amp; 3 until composition of clusters don't change (i.e. convergence)</li> <li>5. Final k clusters</li> </ol>
	Choice of k <ul style="list-style-type: none"> <li>- Min. WSS (Within sum of squares) of clusters</li> <li>- Measures sum of ttl. dist. of (o) in each cluster from their respective centroids</li> </ul>
Decision Tree	
<ul style="list-style-type: none"> <li>- Uses a tree structure to create a rule for classifying data</li> <li>- Sequence of decisions specifying consequences</li> <li>- Used for prediction</li> </ul>	
Build: <ol style="list-style-type: none"> <li>1. Full data with all its var. and (o)</li> </ol> <u>Sequence of if-then statements:</u> <ol style="list-style-type: none"> <li>2. Split data into groups</li> <li>3. Each group split into subgroups</li> <li>4. Splitting con'ts till certain stopping rule is reached or each (o) becomes a group of its own</li> </ol>	
Classification Tree	Regression Tree
<ul style="list-style-type: none"> <li>- o/p var.s usually applied to be categorical in nature</li> <li>- binary decisions like 'yes' or 'no'</li> </ul>	<ul style="list-style-type: none"> <li>- o/p var.s can be numeric or cts</li> </ul>

### Association rules (Unsupervised learning)

- Uncover r/s of form  $X \rightarrow Y$ : when  $X(o)$ , item  $Y(o)$
- Descriptive, not a predictive method
- Itemsets: items bought in the same transaction, webpages visited in the same sitting etc.
- k-itemset: { item 1, item 2, ... , item k }

### Apriori algorithm

- Given collection of items, explores all subsets of items
- Provides subsets which appear more than some pre-defined frequency
- Need transaction database  $D$ , min. length support threshold , max. length an itemset could reach  $N$  (optional)

1. Support	Given itemset $L$ , $\sim$ of $L$ : proportion of transactions that contain $L$ ( $0 < \text{support} < 1$ )
2. Frequent itemset	$\sim$ has items that appear tgt. often enough ( satisfy a min. support criterion, usually set at 0.5 )
3. Downward closure property (Apriori property)	If itemset considered frequent, any subset of freq. itemset also frequent.

Steps:

- Bottom-up approach: Starting from 1-itemset, merges pairs of current itemsets to create new itemsets
- Each stage: Algorithm checks if all itemsets satisfy min. support criterion  
(If not, dropped and not considered in rest of algorithm)
- Process continues till it runs out support or itemsets reach a predefined max. length

### SQL (Structured Query Language)

In-database analytics: describes the processing of data within its repository

- 😊: database can update regularly  $\rightarrow$  produce most up-to-date results
- : Eliminate need to move data from place to place
- : Well-protected/ Security  $\rightarrow$  need to extract to another portal to make changes/ perform analysis

Relational database: part of RDBMS (Relational Database Management system)

- 😊: organizes data in tables with established r/s between tables
- : splitting data into tables  $\rightarrow$  no need to store entire table of info
- : can change a specific part  $\rightarrow$  no need to change everything and store in same place
- 1. Smaller memory to handle them
- 2. Various tables stored on diff. machines
- 3. Changes & corrections easily made
- 4. Relational database: number of duplicates reduced

### Parallel Computation

Distributed:  $\sim$  of computation happens at processor level

- : compilers optimised to distribute computation whenever

Parallel: run a function repeatedly in diff. processors

- : snowFT  $\rightarrow$  run full function in parallel

### Residual bootstrap

Given dataset  $(x_1, x_2, \dots, x_n)$

- Calculate mean, var.
- Resample the data of size  $n$  with replacement (means pts can be repeated)
- Recalc. Data
- Do for a lot of times (thousands?)

Can show that sample distribution will be close to the true distribution of the r.v. looking at