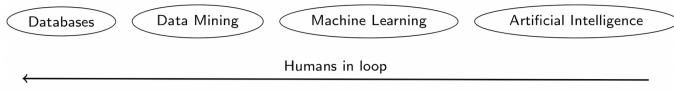


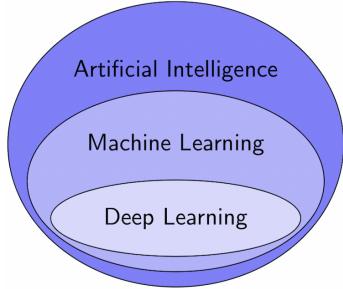
Machine Learning

Introduction

Data Mining vs. Machine Learning



- Data Mining often viewed as closer to databases: automatically extract useful knowledge from large data sets
- Machine Learning often viewed as closer to AI: use computers to automatically detect patterns in data



Fundamentals

- Basic assumptions
 - Training and test data need to be related in some way
 - Training and test data are (approximately) independent and identically distributed (IID)
- Memorization vs learning:
 - You can do well on training data by memorizing it (practice on the exam)
 - You have only learned if you can do well in new situations (the actual exam)
- Parameters vs hyper-parameters
 - We estimate parameters by training a model
 - We tune hyper-parameters using a validation score
- Golden Rule of Machine Learning - the test data cannot influence the training phase in any way, e.g.
 - After tuning hyperparameters on training/validation sets, the selected model is usually retrained on the whole training/validation data, with the test data untouched all the time
 - Imputation on training and test/validation sets are done separately
- No free lunch theorem - there is no best model achieving the best generalization error for every problem

Preprocessing

Data cleaning

- Missing/Incomplete values lacking attribute values, certain attributes, or containing only aggregate data

- Problems
 - Missing completely at random (MCAR)
 - Completely unrelated to the data
 - Potential problem? Small sample size
 - Missing at random (MAR)
 - The fact the data are missing is related not to the missing attribute, but to some other data in the data set
 - Potential problem? Bias due to row-wise deletion
 - Missing not at random (MNAR)
 - There is a reason the data are missing and it is related to the attribute itself
 - Potential problem? Bias due to row-wise deletion

- Solutions (Imputation)
 - Usually ignore the record if the class label is missing
 - A global constant

X	X'
sunny	sunny
cloudy	?
sunny	?
cold	July
...	...
overcast	...
cold	Sat
Sat	June
...	...

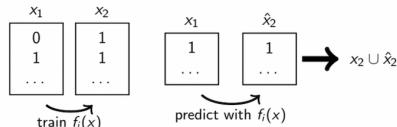
- The attribute mean changes relationship with other variables -> bias in data

X	X'
12	12
2	2
22	22
38	37
...	...
11	11
?	12
?	38
90	90
...	...
2	2
23	23
?	38
?	30
...	...
9	9
11	11
54	54
23	23
...	...

- The attribute mean of the samples belonging to the same class might change relationship with other variables in other classes -> bias in data

$X Y$						$X' Y$					
12	2	22	38	...	1	12	2	22	38	...	1
11	?	?	90	...	0	11	11	54	90	...	0
2	23	?	?	...	1	2	23	22	38	...	1
...
9	11	54	23	...	0	9	11	54	23	...	0

- The most probable value can be inferred by machine learning algorithms



- Matrix factorization is a collaborative filtering method to predict the missing values using the latent features, e.g. singular value decomposition (SVD)

- Decompose the data matrix M such that $M = U\Sigma V^*$
- Create imputed matrix M' by multiplying $U \times \Sigma \times V^*$

- Expectation Maximization

- Use other variables to impute the values (Expectation)
- Check if value is most probable (Maximization)

- Noisy data containing noise, errors, or outliers

- Binning: first sort data and partition into (equal-frequency) bins, then one can be smoothed by bin means, bin median, bin boundaries, etc.

- Clustering

- Inconsistent data containing discrepancies in codes or names

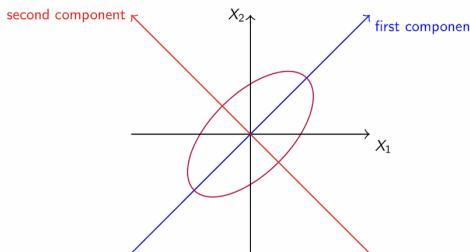
- Intentionally wrong data, e.g. there are a lot of pictures with a GPS location just a bit west of Africa

Data reduction

- Dimensionality reduction - "curse of dimensionality": density and distance between points, which is critical to clustering, outlier analysis, classification, regression becomes less meaningful

- Principal Component Analysis – PCA

- Find a projection that captures the largest amount of variation in data
- The original data are projected onto a much smaller space, resulting in dimensionality reduction
- We find the eigenvectors of the covariance matrix, and these eigenvectors define the new space



- Given N data vectors from n -dimensions, find $k \leq n$ orthogonal vectors (principal components) that can be best used to represent data

- Normalize input data: Each attribute falls within the same range
- Compute k orthonormal (unit) vectors, i.e. principal components
- Each input data (vector) is a linear combination of the k principal component vectors
- The principal components are sorted in order of decreasing "significance" or strength Since the components are sorted, the size of the data can be reduced by eliminating the weak components, i.e. those with low variance (i.e. using the strongest principal components, it is possible to reconstruct a good approximation of the original data)

- Feature selection

- Correlation

- For nominal data, given two attributes A and B with values a_1, \dots, a_c and b_1, \dots, b_r the correlation can be calculated using the χ^2 test:

$$\chi^2 = \sum_{i=1}^c \sum_{j=1}^r \frac{(o_{ij} - e_{ij})^2}{e_{ij}}$$

- o_{ij} is the actual frequency of the event (a_i, b_j)
- e_{ij} is the expected frequency (n is the number of instances)

- Numerical data can be compared using Pearson's correlation coefficient:

$$r_{A,B} = \frac{\sum_{i=1}^n (a_i - \bar{A})(b_i - \bar{B})}{n\sigma_A \sigma_B} = \frac{\sum_{i=1}^n (a_i b_i) - n\bar{A}\bar{B}}{n\sigma_A \sigma_B}$$

where means are \bar{A} and \bar{B} , number of instances is n , and standard deviations are σ_A and σ_B

- Relief is a feature selection algorithm that rewards the differences (different class labels) and penalize the similarities (the same class label)

Input: Data set with N_f attributes and N_i instances that belong to one of two classes, and parameter $N_r < N_i$

First normalize the data

Create a weight vector W with one weight $w_i \in W$ for each attribute

Initialize the weights to 0

for $j \in 1 \dots N_f$ do

 Randomly select instance $X = [x_1, \dots, x_n]$

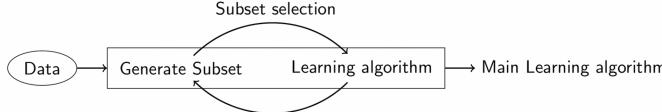
 Choose instance $H = [h_1, \dots, h_n]$ as the closest neighbour of X in the same class ($nearHit$)

```

    choose instance  $v = [v_1, \dots, v_m]$  as the closest neighbour of  $X$  in the same class (nearest)
    choose instance  $M = [m_1, \dots, m_n]$  as the closest neighbour of  $X$  in the other class (nearMiss)
    for  $i \in 1 \dots N_d$  do
        |  $w_i = w_i - (x_i - h_i)^2 + (x_i - m_i)^2$ 
    end
    end
    for  $i \in 1 \dots N_d$  do
        |  $w_i = \frac{w_i}{N_r}$ 
    end

```

- Wrappers



- generate a subset of the features and evaluate the performance of the classifier on the subset
- Add or remove attributes from the subset and see if the performance of the classifier improves
- Numerosity reduction
 - Regression and Log-Linear Models
 - Histograms, clustering, sampling
 - Data cube aggregation Data compression
- Data compression

Transformation and discretization

- Normalization
 - Min-max normalization to the range of (new_minA, new_maxA)
 - $v' = \frac{v - min_A}{max_A - min_A} (new_max_A - new_min_A) + new_min_A$
 - Z-score normalization – mean μ , standard deviation σ
 - $v' = \frac{v - \mu_A}{\sigma_A}$
 - Normalization by decimal scaling
- Where j is the smallest integer such that $\text{Max}(|v'|) < 1$
- Imbalanced data
 - Randomly under-sample the majority class
 - Randomly over-sample the minority class
- Hierarchy generation
 - Problems
 - Nominal – values from an unordered set, e.g. colour
 - Ordinal – values from an ordered set, e.g. rank
 - Numeric – real numbers, e.g. integers or reals
 - Solutions
 - Equal-width (distance) partitioning, sensitive to outliers and skewed data
 - Equal-depth (frequency) partitioning, tricky handling categorical data

Classification

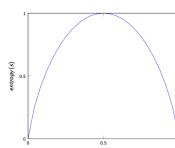
Decision Tree can be seen a nested sequence of “if-else” decisions based on the features (splitting rules) and a class label as a return value at the end of each sequence

- Train a decision stump - a simple decision tree with one splitting rule based on thresholding one feature
 - Define a 'score' for the rule, e.g. "Information Gain" is the most common score in practice to decrease the entropy (the degree of disorder or randomness in the system)

$$\text{information gain} = \frac{\text{entropy}(y)}{\text{entropy before split}} - \frac{n_{\text{yes}}}{n} \underbrace{\text{entropy}(y_{\text{yes}})}_{\text{entropy examples satisfying rule}} - \frac{n_{\text{no}}}{n} \underbrace{\text{entropy}(y_{\text{no}})}_{\text{entropy examples satisfying rule}}$$

with

$$\text{entropy}(s) = -p_{\text{pos}} \log_2 p_{\text{pos}} - p_{\text{neg}} \log_2 p_{\text{neg}}$$



- Search for the rule with the best scores
- Greedy Recursive Splitting
 - Fit a decision stump to each leaf's data
 - Add these stumps to the tree
 - Stop when
 - Clean split (leaves) with only one class label
 - User-defined maximum depth
 - Prune when the information gain is low for several levels and then becomes high again

Reduced Error Pruning

Input: decision Tree T ; labelled data D

Output: Pruned tree T'

for every internal node N of T , starting from the bottom do

```

     $T_N \leftarrow \text{subtree of } T \text{ rooted at } N;$ 
     $D_N \leftarrow \{x \in D | x \text{ is covered by } N\};$ 
    if accuracy of  $T_N$  over  $D_N$  is worse than majority class in  $D_N$  then
        | replace  $T_N$  in  $T$  by a leaf labelled with the majority class in  $D_N$ ;

```

```

    | end
end
return pruned version T

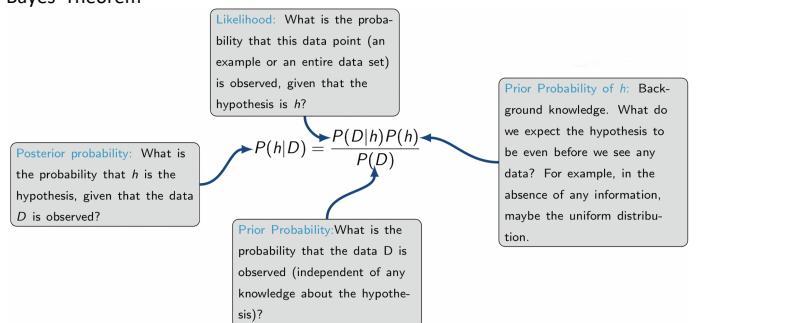
```

Ensemble

- Definition - an ensemble of classifiers is a set of classifiers whose individual decisions are combined in some way to classify new examples
 - The individual classifiers disagree with one another
 - The individual classifiers make uncorrelated errors at rates less than random
 - Bootstrapping to reduce the variance of an individual classifier
 - Construct ensembles
 - Manipulating the training set for unstable learning algorithms
 - Resample with replacement (bagging)
 - Resample without replacement (Cross-validated committees)
 - Resample with adjusted weights (boosting -> overfitting)
 - Resample with penalization (gradient to take steps in the opposite direction -> underfitting)
 - Manipulating the input features for correlated input features, e.g. only randomly choose m out of p features (random forest, e.g. $m \approx \sqrt{p}$)
 - Manipulating the output features, e.g. Error Correcting Output Coding (ECOC) for using binary classification models on multi-class classification tasks
- | Class | 1 | 2 | 3 | 4 | 5 | 6 | 7 | 8 | 9 | 10 | 11 | 12 | 13 | 15 |
|-------|---|---|---|---|---|---|---|---|---|----|----|----|----|----|
| A | 0 | 0 | 0 | 1 | 0 | 0 | 1 | 0 | 1 | 1 | 0 | 1 | 1 | 1 |
| B | 0 | 0 | 1 | 0 | 0 | 1 | 0 | 1 | 0 | 1 | 1 | 0 | 1 | 1 |
| C | 0 | 1 | 0 | 0 | 1 | 0 | 0 | 1 | 1 | 0 | 1 | 1 | 0 | 1 |
| D | 1 | 0 | 0 | 0 | 1 | 1 | 1 | 0 | 0 | 0 | 1 | 1 | 1 | 0 |
- Injecting randomness
 - Into decision tree split criteria so they chooses randomly among the best n tests at each node
 - Into ANN bootstrap sampling of training data to add Gaussian noise to the input features
 - Combine classifiers
 - Unweighted
 - Weighted
 - Stacking - use outputs of n classifiers as attributes for target
 - Gating = stacking + weighted/unweighted
 - Random forest
 1. Sample N (size of training set) cases at random - with replacement, from the original data. This sample will be the training set for growing the tree.
 2. For M input variables, a number m < M is specified such that at each node, m variables are selected at random out of the M and the best split on these m is used to split the node. The value of m is held constant during the forest growing.
 3. Each tree is grown to the largest extent possible. There is no pruning
 - Comparisons
 - Random Forest
 - Bagging and Random Forests can be easily run in parallel, but not for Boosting for there is no random sample drawn.
 - No need for cross-validation or a separate test set to get an unbiased estimate of the test set error for as it can be tested again the out of bag (OOB) samples during the run
 - A good m needs to be tuned to balance trade-offs between correlation and strength of each tree in the forest
 - Variable importance can be ranked by randomly permuting values of m variables in the OOB cases, and comparing them to the untouched OOB cases (only work if variables are independent)
 - AdaBoost: natively reweighted, so it's sensitive to class errors
 - Gradient Boost: can be penalised by L1 or L2 regularisation
 - XG Boost: improve the performance of Gradient Boost series

Bayesian Learning

Bayes' Theorem



- D: The event that we observed this particular data set
- h: The event that the hypothesis h is the true hypothesis

Applications

- Explicit manipulation of probabilities, e.g. among the most practical approaches to certain types of learning problems

- Useful framework for understanding learning methods that do not explicitly manipulate probabilities, e.g. determine conditions under which algorithms output the most probable hypothesis: justification of the error functions in ANNs, justification of the inductive bias of decision trees, etc

Maximum A Posteriori (MAP)

In many learning scenarios, the learner considers some set of candidate hypotheses H and is interested in finding the most probable hypothesis $h \in H$ given the observed training data D

- The maximum a posteriori (MAP) hypothesis h_{MAP} where $P(D)$ is dropped because it is a constant independent of h

$$\begin{aligned} h_{MAP} &= \arg \max_{h \in H} P(h|D) \\ &= \arg \max_{h \in H} \frac{P(D|h)P(h)}{P(D)} \\ &= \arg \max_{h \in H} P(D|h)P(h) \end{aligned}$$

- Example

- Consider a medical diagnosis problem in which there are two alternative hypotheses
 - The patient has a particular form of cancer (denoted by *cancer*)
 - The patient does not have a cancer (denoted by \neg *cancer*)
- The available data is from a particular laboratory with two possible outcomes: \oplus (positive) and \ominus (negative)

$$\begin{array}{ll} P(\text{cancer}) = 0.008 & P(\neg\text{cancer}) = 0.992 \\ P(\oplus|\text{cancer}) = 0.98 & P(\ominus|\text{cancer}) = 0.02 \\ P(\oplus|\neg\text{cancer}) = 0.03 & P(\ominus|\neg\text{cancer}) = 0.97 \end{array}$$

- Suppose a new patient is observed for whom the lab test returns a positive (\oplus) result
- Should we diagnose the patient as having cancer or not?

$$\begin{aligned} P(\text{cancer}|\oplus) &\sim P(\oplus|\text{cancer})P(\text{cancer}) & = 0.98 * 0.008 = 0.0078 \\ P(\neg\text{cancer}|\oplus) &\sim P(\oplus|\neg\text{cancer})P(\neg\text{cancer}) & = 0.03 * 0.992 = 0.0298 \\ \Rightarrow h_{MAP} &= \neg\text{cancer} \end{aligned}$$

- The exact posterior probabilities can be determined by normalizing the above probabilities to 1

$$\begin{aligned} P(\text{cancer}|\oplus) &= \frac{0.0078}{0.0078 + 0.0298} = 0.21 \\ P(\neg\text{cancer}|\oplus) &= \frac{0.0298}{0.0078 + 0.0298} = 0.79 \end{aligned}$$

Maximum Likelihood (ML)

- The equation above can be further simplified assuming that every hypothesis is equally probable a priori. This is called maximum likelihood (ML) hypothesis h_{ML}

$$h_{ML} = \arg \max_{h \in H} P(D|h)$$

- Least-Squared Error

- Problem: learning continuous-valued target functions (e.g. neural networks, linear regression, etc.). Under certain assumptions any learning algorithm that minimizes the squared error between the output hypothesis and the training data, will output a ML hypothesis
 - $(\forall h \in H)[h : X \rightarrow \mathbb{R}]$ and training examples of the form $< x_i, d_i >$
 - unknown target function $f : X \rightarrow \mathbb{R}$
 - m training examples, where the target value of each example is corrupted by random noise drawn according to a Normal probability distribution with zero mean $(d_i = f(x_i) + e_i)$

- Solution:
 - It is common to maximize the less complicated logarithm, which is justified because of the monotonicity of this function

$$h_{ML} = \arg \max_{h \in H} \sum_{i=1}^m \log \frac{1}{\sigma \sqrt{2\pi\sigma^2}} - \frac{(d_i - h(x_i))^2}{2\sigma^2}$$

- The first term in this expression is a constant independent of h and can therefore be discarded.

$$h_{ML} = \arg \max_{h \in H} \sum_{i=1}^m -\frac{(d_i - h(x_i))^2}{2\sigma^2}$$

- Maximizing this negative term is equivalent to minimizing the corresponding positive term.

$$h_{ML} = \arg \min_{h \in H} \sum_{i=1}^m \frac{(d_i - h(x_i))^2}{2\sigma^2}$$

- Finally, all constants independent of h can be discarded.

$$h_{ML} = \arg \min_{h \in H} \sum_{i=1}^m (d_i - h(x_i))^2$$

\Rightarrow the h_{ML} is one that minimizes the sum of the squared errors

Minimum Description Length (MDL)

- Occam's razor: choose the shortest explanation for the observed data
- Interpret the definition of h_{MAP} in the light of information theory concepts, we can see that short hypotheses are preferred, assuming a particular representation scheme for encoding hypotheses and data

$$\begin{aligned} h_{MAP} &= \arg \max_{h \in H} P(D|h)P(h) \\ &= \arg \max_{h \in H} \log_2 P(D|h) + \log_2 P(h) \end{aligned}$$

$$= \arg \min_{h \in H} -\log_2 P(D|h) - \log_2 P(h)$$

- A basic result of information theory
 - Problem: design a code C to transmit messages drawn at random
 - Probability of encountering message i is p_i
 - Interested in the most compact code C
 - Shannon and Weaver (1949) showed that the optimal code assigns $-\log_2 p_i$ bits to encode message i
 - $L_C(i) \approx$ description length of message i with respect to C
 - Solution: Minimum description length principle

$$h_{MAP} = \arg \min_{h \in H} L_{C_H}(h) + L_{C_{D|h}}(D|h)$$

- $L_{CH}(h) = -\log_2 P(h)$, where CH is the optimal code for hypothesis space H
- $L_{CD|h}(D|h) = -\log_2 P(D|h)$, where CD|h is the optimal code for describing data D assuming that both the sender and receiver know hypothesis h
- MDL principle provides a way for trading off hypothesis complexity for the number of errors committed by the hypothesis, e.g. decision tree learning
 - C_H might be some obvious encoding, in which the description length grows with the number of nodes and with the number of edges
 - Choice of $C_{D|h}$?
 - Sequence of instances $\langle x_1, \dots, x_m \rangle$ is known to the transmitter and the receiver
 - We need only to transmit the classifications $\langle f(x_1), \dots, f(x_m) \rangle$
 - If h correctly predicts the classification, no transmission is necessary ($L_{C_{D|h}}(D|h) = 0$)
 - In case of missclassified examples, for each missclassification a message has to be sent that identifies this example (at most $\log_2 m$ bits) as well as its correct classification (at most $\log k$ bits, where k is the number of possible classifications)

Bayes Optimal Classifier

- This is different from the MAP framework that seeks the most probable hypothesis (model). Instead, we are interested in making a specific prediction, so the most probable classification is obtained by combining the predictions of all hypotheses, weighted by their posterior probabilities

$$\arg \max_{v_j \in V} P(v_j|D) = \arg \max_{v_j \in V} \sum_{h_i \in H} P(v_j|h_i)P(h_i|D)$$

where $P(v_j|D)$ is the probability that the correct classification is v_j

- Example

$$V = \{\oplus, \ominus\}$$

$$P(h_1|D) = 0.4, P(\ominus|h_1) = 0, P(\oplus|h_1) = 1$$

$$P(h_2|D) = 0.3, P(\ominus|h_2) = 1, P(\oplus|h_2) = 0$$

$$P(h_3|D) = 0.3, P(\ominus|h_3) = 1, P(\oplus|h_3) = 0$$

therefore

$$P(\oplus|D) = \sum_{h_i \in H} P(\oplus|h_i)P(h_i|D) = 0.4$$

$$P(\ominus|D) = \sum_{h_i \in H} P(\ominus|h_i)P(h_i|D) = 0.6$$

and

$$\arg \max_{v_j \in \{\oplus, \ominus\}} \sum_{h_i \in H} P(v_j|h_i)P(h_i|D) = \ominus$$

Naïve Bayes Classifier

- The computational cost of Bayes optimal classifier is to weigh in the predictions of all hypotheses. In the case where each instance x is described by a conjunction of attribute values $\langle a_1, a_2, \dots, a_n \rangle$ and where the target function $f(x)$ can take on any value from some finite set V , this can be simplified assuming attribute values are conditionally independent given the target value:

$$v_{NB} = \arg \max_{v_j \in V} P(v_j) \prod_i P(a_i|v_j)$$

- $P(v_j)$ can be estimated by counting the frequency of v_j in D
- The number of terms is $|$ distinct attribute values $| \times |V| + |V|$
- Document Classification
 - For documents presented by features x_1, x_2, \dots, x_n , let X be the set of unique feature values (e.g. unique words)

$$c_{MAP} = \arg \max_{c \in C} P(c) \prod_i P(x_i|c)$$

where document features can be based on word counts

$$P(x_i|c) = \frac{count(x_i, c)}{\sum_{x_j \in X} count(x_j, c)}$$

- Laplace smoothing for unknown words, i.e. $count(x_i, c) = 0$ for all classes, a constant is added

$$P(x_i|c) = \frac{count(x_i, c) + 1}{\sum_{x_j \in X} count(x_j, c) + 1}$$

$$\sum_{x_j \in X} (count(x_j, c) + 1)$$

- Example

	Doc.ID	Words	class	$P(c) = \frac{N_c}{N}$
Training	d1	Kiwi, Sheep, Kiwi	NZ	
	d2	Kiwi, Kiwi, Bird	NZ	
	d3	Kiwi, Auckland	NZ	
	d4	Munich, Oktoberfest, Kiwi	DE	
Test	d5	Kiwi, Kiwi, Munich, Oktoberfest	NZ	

- Priors

$$P(NZ) = 3/4$$

$$P(DE) = 1/4$$

- Choosing a class for d5

$$P(NZ|d5) \sim 3/4 * (3/7)^3 * 1/14 * 1/14 \approx 0.0003$$

$$P(DE|d5) \sim 1/4 * (2/9)^3 * 2/9 * 2/9 \approx 0.0001$$

- Conditional Probabilities

$$P(Kiwi|NZ) = (5+1)/(8+6) = 3/7$$

$$P(Munich|NZ) = (0+1)/(8+6) = 1/14$$

$$P(Oktoberfest|NZ) = (0+1)/(8+6) = 1/14$$

$$P(Kiwi|DE) = (1+1)/(3+6) = 2/9$$

$$P(Munich|DE) = (1+1)/(3+6) = 2/9$$

$$P(Oktoberfest|DE) = (1+1)/(3+6) = 2/9$$

Instance-based Learning

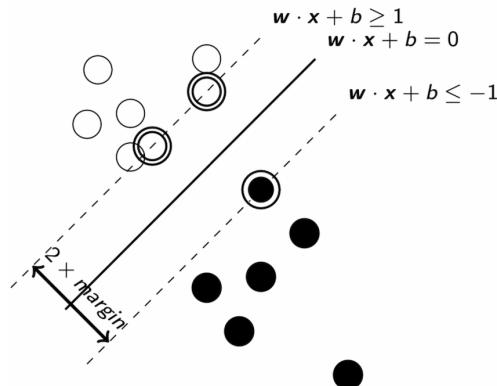
A class of algorithms that do not learn a function but rather directly compare to known examples

K-Nearest Neighbour (KNN)

- Compute the pair-wise distances (Euclidean, Mahalanobis distance, etc) to the labelled examples, and find the nearest k neighbours
- Choose the class label of the majority votes from the k nearest neighbours

Support Vector Machines (SVM)

- Hard margin: maximum margin hyperplane achieves maximal separation between the classes where the encircled training instances define the margin borders for each class -> maximal margin classifier



- Training data $\{(x_k = [x_{k1}, \dots, x_{km}], y_k)\} y_k \in \{-1, +1\}, k \in [1, n]$
- Classify instance x_k as $+1$ if $w \cdot x_k + b \geq 1$, -1 if $w \cdot x_k + b \leq -1$, where $w \cdot x_k = \sum_{i=1}^m w_i x_{ki}$
- The distance of point x to a hyperplane defined with parameters w and b is $\frac{|w \cdot x + b|}{\sqrt{w \cdot w}}$ where the distance of $x - x_+$ to the decision boundary is reduced $\frac{1}{\sqrt{w \cdot w}}$
- Soft margin allows misclassification errors -> support vector classifier/ soft margin classifier

It maximizes using Lagrange multiplier

$$L(\alpha) = \sum_{k=1}^n \alpha_k - \frac{1}{2} \sum_{k=1}^n \sum_{q=1}^n \alpha_k \alpha_q y_k y_q x_k \cdot x_q$$

with respect to α subject to the constraints $0 \leq \alpha \leq C$

- Complexity parameter $C > 0$ trades off training error for model complexity (number of support vectors, i.e. the observations that violate the margin)
- $C \rightarrow \infty$ we get the maximal margin classifier
- Non-linear kernel transforms the dot product in the decision making process -> support vector machine

$$d(\mathbf{x}) = b + \sum_{i=1}^l \alpha_i y_i \mathbf{x}_i \cdot \mathbf{x}$$

becomes

$$d(\mathbf{x}) = b + \sum_{i=1}^l \alpha_i y_i (\mathbf{x}_i \cdot \mathbf{x})^h$$

- Number of support vectors l
- Parameters α_i , b (estimated by training / optimization)
- Test instance vector \mathbf{x}
- Support vector \mathbf{x}_i
- An instance can be classified by using $\hat{y}(x) = sign(d(x))$
- Different kernels to encode expert knowledge
 - Polynomial kernels: $K(\mathbf{x}_i, \mathbf{x}) = (\gamma \mathbf{x}_i^T \mathbf{x} + 1)^d$, $\gamma > 0$, $d = 1, 2, \dots$ -> image processing
 - Radial kernels: $K(\mathbf{x}_i, \mathbf{x}) = \exp(-\gamma \|\mathbf{x}_i - \mathbf{x}\|^2)$, $\gamma > 0$ -> general purpose
 - Sigmoid kernels: $K(\mathbf{x}_i, \mathbf{x}) = \tanh(\gamma \mathbf{x}_i^T \mathbf{x} + 1)$, $\gamma > 0$ where $\tanh(x) = (e^{2x} - 1)/(e^{2x} + 1)$ -> neural network

Evolutionary Algorithms

Framework

1. Random generation of an initial population
2. Reckoning of a fitness value for each subject. It will directly depend on the distance to the optimum.
3. Reproduction of the population based on fitness values.
4. If requirements are met, then stop. Otherwise go back to

Generic Algorithm (GA)

- Biological foundations
 - Phenotype is determined by genotype and learning/interaction with the environment
 - Genotype is the collection of genes of an organism.
 - Phenotype is the morphology, physiology and behaviour of an organism.
 - Evolution is a search process for phenotypes
 - Chromosomes (DNA molecules) can be interpreted as character strings in nature's base-4 alphabet - nucleotides: A (adenine), C (cytosine), G (guanine) and T (thymine)
 - The evolutionary process performs natural selection and genetic operations (mutation & crossover)
- Genetic Algorithms are highly parallel search algorithms inspired by evolutionary processes

$GA(Fitness, Fitness_threshold, p, r, m)$

Fitness: A function that assigns an evaluation score, given a hypothesis.

Fitness_threshold: A threshold specifying the termination criterion.

p: The number of hypotheses to be included in the population.

r: The fraction of the population to be replaced by Crossover at each step.

m: The mutation rate.

• Initialize population: $P \leftarrow$ Generate p hypotheses at random

• Evaluate: For each h in P , compute $Fitness(h)$

• While $\max_h Fitness(h) < Fitness_threshold$ do

 Create a new generation, P_S :

 1. Select: Probabilistically select $(1-r)p$ members of P to add to P_S . The probability $Pr(h_i)$ of selecting hypothesis h_i from P is given by

$$Pr(h_i) = \frac{Fitness(h_i)}{\sum_{j=1}^p Fitness(h_j)}$$

 2. Crossover: Probabilistically select $\frac{r}{2}p$ pairs of hypotheses from P , according to $Pr(h_i)$ given above. For each pair, (h_1, h_2) , produce two offspring by applying the Crossover operator. Add all offspring to P_S .

 3. Mutate: Choose m percent of the members of P_S with uniform probability. For each, invert one randomly selected bit in its representation.

 4. Update: $P \leftarrow P_S$.

 5. Evaluate: for each h in P , compute $Fitness(h)$

• Return the hypothesis from P that has the highest fitness.

◦ Population models depending on the proportion of the population replaced (generation gap)

- Generational Model (GGA) - Each individual survives 1 generation (1.0)
- Steady State Model (SSGA) - One offspring generated per generation with one member replaced (1/pop_size)

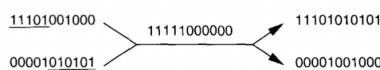
◦ Fitness function

- Accuracy of the hypothesis over the training data
- Number of games won by the individual when playing against other individuals in the current population

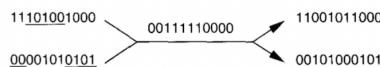
◦ Generic Operators

Initial strings *Crossover Mask* *Offspring*

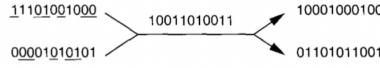
Single-point crossover:



Two-point crossover:



Uniform crossover:



Point mutation:



- Applications - balancing the trade-off of changing variables, e.g. self-driving cars
 - Children of individuals with high fitness will have high fitness
 - Blocks of close genes - assuming uniform crossover is not used, the closer two genes are to each other, the more likely they are to be passed together to the same child
- Limitation - early convergence to a local minima due to the overcrowding of the fittest gene
 - Sharing (fitness sharing causes individuals that are similar to many others to have their fitness reduced, making them less likely to be selected) for GA and crowding (the current population is sampled to find an individual that is "close" to the new offspring, then this closest individual is then replaced) for SSGA

- Tournament Selection - pick k members at random, then select the best of these, can also use p to determine whether the fittest gene wins

Swarm Intelligence

- Biological foundations - collective behaviours result from the local interactions of the individuals with each other and/or with their environment, e.g. colonies of ants and termites, schools of fish, flocks of birds, bacterial growth, herds of land animals.
- Ant Colony Optimisation (ACO) - initially, ants wander randomly and return to their colony upon finding food, laying down pheromone trails. The pheromone trails get reinforced if any other ant follows the trail, finds food and returns to the colony, so a short path gets marched over faster with the pheromone density remains high as it is laid on the path as fast as it can evaporate -> real-time adaptability, e.g. network routing, transportation systems, etc.
 - The choice of a solution component from $N(s^p)$ is done probabilistically at each construction step, e.g.

$$p(c_{ij} | s^p) = \frac{\tau_{ij}^\alpha \cdot \eta_{ij}^\beta}{\sum_{c_{il} \in N(s^p)} \tau_{il}^\alpha \cdot \eta_{il}^\beta}, \forall c_{ij} \in N(s^p)$$

where τ_{ij} and η_{ij} are the pheromone value and the heuristic value associated with the component.

α and β are positive real parameters whose values determine the relative importance of pheromone versus heuristic information.

- Pheromone update
 - Decrease all the pheromone values through pheromone evaporation associated with bad solutions
 - Increase the pheromone levels associated with a chosen set of good solutions

$$T_{ij} \leftarrow (1 - \rho) \cdot T_{ij} + \sum_{k=0}^M (\Delta T_{xy}^k)$$

where $\rho \in (0, 1]$ is a parameter called evaporation rate, and ΔT_{xy}^k is the amount of pheromone left by the kth ant on the edge between x and y

- Particle Swarm Optimization (PSO) - bird flocking or fish schooling
 - A problem is given, and some way to evaluate a proposed solution to it exists in the form of a fitness function.
 - A communication structure or social network is also defined, assigning neighbours for each individual to interact with.
 - Then a population of individuals defined as random guesses at the problem solutions is initialized – candidate solutions.
 - The swarm is typically modelled by particles in multidimensional space that have a position and a velocity

$$v_{i,j} \leftarrow c_0 v_{i,j} +$$

$$c_1 r_1 (globalbest_j - x_{i,j}) +$$

$$c_2 r_2 (localbest_{i,j} - x_{i,j}) +$$

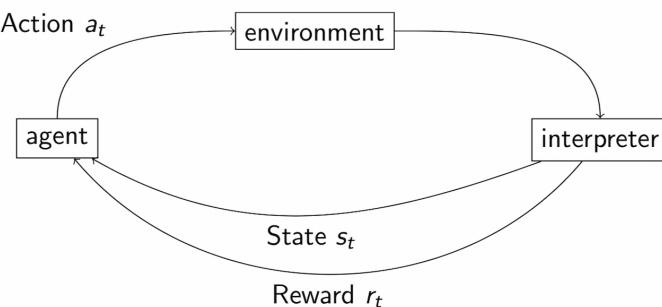
$$c_3 r_3 (neighborhoodbest_j - x_{i,j})$$

$$x_{i,j} \leftarrow x_{i,j} + v_{i,j}$$

i is the particle and j is the dimension

Reinforcement Learning

to choose actions that maximize the total rewards -> a sequential decision making



Markov Decision Process

- Mathematical formulation
 - S Set of possible states $s_t \in S$
 - A Set of possible actions $a_t \in A$
 - R Distribution of the current reward given (state, action) pair
 - P Distribution of the next state given (state, action) pair
 - γ Discount factor with $0 \leq \gamma < 1$ weighting the impact of future rewards
- Assumptions
 - The next state is only related to the current state and action:

$$P(s_{t+1} | s_t, a_t, s_{t-1}, a_{t-1}, \dots) = P(s_{t+1} | s_t, a_t)$$
 - The reward is only related to current state and action:

$$R(r_t | s_t, a_t, s_{t-1}, a_{t-1}, \dots) = R(r_t | s_t, a_t)$$

- The task is learning a policy $\pi : S \rightarrow A$ for choosing actions that maximizes

$$E[r + \gamma r_{t+1} + \gamma^2 r_{t+2} + \dots] \text{ for every possible starting state } s:$$

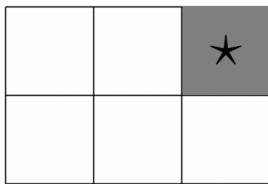
- Value function

$$V^\pi(s) = \mathbb{E} \left[\sum_{t \geq 0} \gamma^t r_t \right]$$

- Optimal policy

$$\pi^* = \arg \max_{\pi} V^\pi(s), \forall s \in S$$

- Example: grid world where one reaches the terminal state (greyed out) in least number of actions.



1. Rewards

0	0	100
0	0	100
0	0	0

2. Possible actions

↔↔	↔↔	*
↔↔	↔↔	↑↔

3. Optimal policy

→	→	*
↑→	↑→	↑

4. Suppose $\gamma = 0.9$ and optimal policy

90	0	0	100	100	*
0			0		
81	0	0	90	0	100

Q-Learning

- Value iteration for $V^{\pi^*}(s)$ (can be abbreviated as $V^*(s)$)

■ Assuming $P(S_{t+1}|S_t, A)$ is known

Initialize $V(s)$ arbitrarily

while Not Done **do**

```

foreach  $s \in S$  do
    foreach  $a \in A$  do
         $| Q(s, a) \leftarrow r(s, a) + \gamma \sum_{s' \in S} P(s'|s, a) V(s')$ 
    end
     $| V(s) \leftarrow \max_a Q(s, a)$ 
end

```

end

■ $V(s)$ converges to $V^*(s)$

- Interestingly, value iteration works even if we randomly traverse the environment instead of looping through each state and action methodically, but we must still visit each state infinitely often on an infinite run -> online learning as agent randomly roams
- If max (over states) difference between two successive value function estimates is less than ϵ , then the value of the greedy policy differs from the optimal policy by no more than $2\epsilon\gamma/(1-\gamma)$

$$\max_s |V_{i+1}(s) - V_i(s)| < \epsilon \Rightarrow \max_s |V^*(s) - V_i(s)| < \frac{2\epsilon\gamma}{1-\gamma}$$

- Q-learning when $P(S_{t+1}|S_t, A_t)$ is unknown

■ ■ ■

```

toreach  $s, a$  do
    | initialize table entry  $\hat{Q}(s, a) \leftarrow 0$ 
end
observe current state  $s$ 
while forever do
    | select action  $a$  and execute it
    | receive immediate reward  $r$ 
    | observe new state  $s'$ 
    | update the table entry for  $\hat{Q}(s, a)$ :
    |  $\hat{Q}(s, a) \leftarrow r + \gamma \max_{a'} \hat{Q}(s', a')$ 
    |  $s \leftarrow s'$ 
end

```

- Q and V^* are closely related:

$$V^*(s) = \max_{a'} Q(s, a')$$

- \hat{Q} is to take into account the maximum reward for the next state after executing a action with the probability that's proportional to rewards (so there's still a chance to jump round to find the global optimum instead of being trapped of local optima) to approximate Q :

$$\hat{Q}(s, a) \leftarrow \mathbb{E}[r(s, a)] + \gamma \max_{a'} \hat{Q}(s', a')$$

- \hat{Q} will always converge to the true Q function
 - i. The system is a deterministic MDP
 - ii. The immediate reward values are bounded: $\forall s, \forall a, \exists c > 0 : |r(s, a)| < c$
 - iii. Agents select actions such that every state-action pair is visited infinitely often
- The training rule for the nondeterministic Q -learning

Q Learning generalizes to nondeterministic worlds

Alter training rule to

$$\hat{Q}_n(s, a) \leftarrow (1 - \alpha_n) \hat{Q}_{n-1}(s, a) + \alpha_n [r + \gamma \max_{a'} \hat{Q}_{n-1}(s', a')]$$

where

$$\alpha_n = \frac{1}{1 + \text{visits}_n(s, a)},$$

- If α_n is set to 1, same rule as the deterministic case
- We can still prove convergence of \hat{Q} to Q [Watkins and Dayan, 1992]

Association Analysis

Association Rule Mining

Given a set of transactions, find rules that will predict the occurrence of an item based on the occurrences of other items in the transaction.

- Definitions
 - Frequent Itemset: an itemset whose support is greater than or equal to a minsup threshold
 - k-itemset: an itemset that contains k items
 - Support count (σ): frequency of occurrence of an itemset
 - Support: fraction of transactions that contain an itemset
 - Association Rule: a rule whose confidence is greater than or equal to minconf threshold
 - An implication expression of the form $X \rightarrow Y$, where X and Y are itemsets
 - Rule Evaluation Metrics
 - Support (s): fraction of transactions that contain both X and Y
 - Confidence (c): measures how often items in Y appear in transactions that contain X
 - Example

TID	Items
1	Bread, Milk
2	Bread, Diaper, Beer, Eggs
3	Milk, Diaper, Beer, Coke
4	Bread, Milk, Diaper, Beer
5	Bread, Milk, Diaper, Coke

Example: $\{\text{Milk, Diaper}\} \Rightarrow \{\text{Beer}\}$

$$s = \frac{\sigma(\text{Milk, Diaper, Beer})}{|T|} = \frac{2}{5} = 0.4$$

$$c = \frac{\sigma(\text{Milk, Diaper, Beer})}{\sigma(\text{Milk, Diaper})} = \frac{2}{3} = 0.67$$

- Mining Tasks
 - Frequent Itemset Generation, given d unique items, the total number of itemsets M :

$$M = 2^d$$
 - Rule Generation, given d unique items, the total number of possible association rules R :
- $$R = \sum_{k=1}^{d-1} \left[\binom{d}{k} \times \sum_{j=1}^{d-k} \binom{d-k}{j} \right]$$
- $$= 3^d - 2^{d+1} + 1$$
- Given a frequent itemset L , find all non-empty subsets $f \subset L$ such that $f \rightarrow L - f$ satisfies the minimum confidence requirement.
 - If $|L| = k$, then there are $2k - 2$ candidate association rules (ignoring $L \rightarrow \emptyset$ and $\emptyset \rightarrow L$)

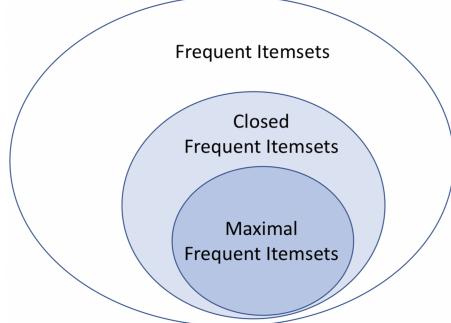
Apriori Algorithm for Frequent Itemset Generation

- Apriori principle - If an itemset is frequent, then all of its subsets must also be frequent:

$$\forall X, Y : (X \subseteq Y) \Rightarrow s(X) \geq s(Y)$$

- Support of an itemset never exceeds the support of its subsets
- This is known as the anti-monotone property of support

- Maximal Frequent vs Closed Frequent



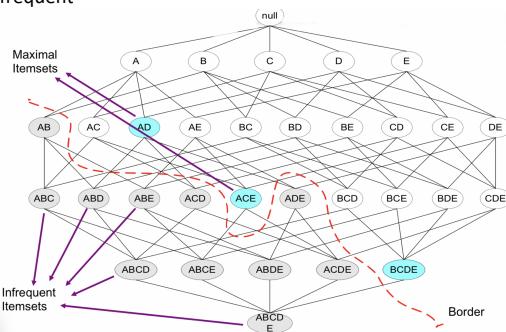
- Closed: an itemset X is closed if none of its immediate supersets has the same support as the itemset X .

TID	Items
1	{A,B}
2	{B,C,D}
3	{A,B,C,D}
4	{A,B,D}
5	{A,B,C,D}

Itemset	Support
{A}	4
{B}	5
{C}	3
{D}	4
{A,B}	4
{A,C}	2
{A,D}	3
{B,C}	3
{B,D}	4
{C,D}	3

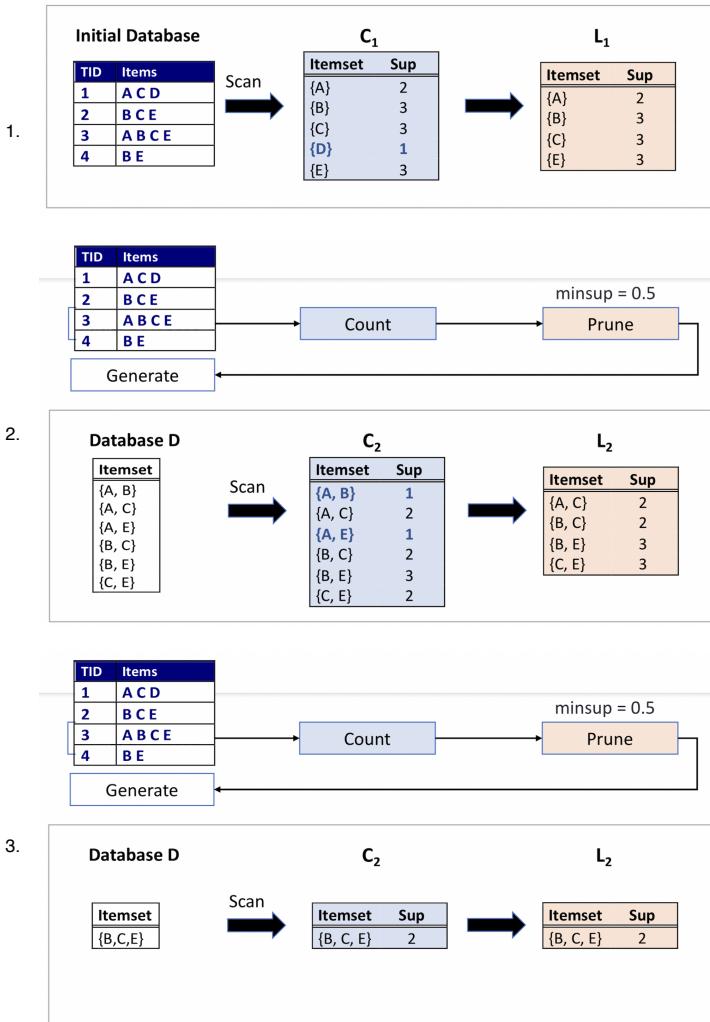
Itemset	Support
{A,B,C}	2
{A,B,D}	3
{A,C,D}	2
{B,C,D}	2
{A,B,C,D}	2

- Maximal: an itemset is maximal frequent if none of its immediate supersets is frequent



- Algorithm
 - Let $k=1$
 - Generate frequent itemset sof length 1
 - Repeat until no new frequent itemsets are identified
 - Generate length $(k+1)$ candidate itemsets from length k frequent itemsets
 - Prune candidate itemsets containing subsets of length k that are infrequent
 - Count the support of each candidate by scanning the DB
 - Eliminate candidates that are infrequent, leaving only those that are frequent
 - Example (automatically):





- Example (manually)

	A	B	C	D	E	F	G	H	I	J
Transactions	1									
2		■			■	■				■
3			■					■		
4				■					■	
5					■					
6						■				
7							■			
8										
9								■		
10										

Support threshold (by count) : 5

- Frequent itemsets: {F}
- Maximal itemsets: {F} .

Support threshold (by count) : 4

- Frequent itemsets: {E}, {F}, {E,F}, {I}
- Maximal itemsets: {E,F}, {I}

Support threshold (by count) : 3

- Frequent itemsets: All subsets of {C,D,E,F} + {J}
- Maximal itemsets: {C,D,E,F}, {J}

Pattern Evaluation for Rule Generation

- Confidence($X \rightarrow Y$) > support(Y), otherwise, rule will be misleading because having item X reduces the chance of having item Y in the same transaction

	Coffee	Coffee	
Tea	15	5	20
Tea	75	5	80
	90	10	100

Association Rule: Tea → Coffee

$$\text{Confidence} = P(\text{Coffee} | \text{Tea}) = 15/20 = 0.75$$

but $P(\text{Coffee}) = 0.9$, which means knowing that a person drinks tea reduces the probability that the person drinks coffee!

- Confidence of rules generated from the same itemset has an anti-monotone property, i.e. confidence is anti-monotone with reference to (w.r.t) items on the right hand side of the rule

e.g., $L = \{A, B, C, D\}$,

$$\text{conf}(ABC \rightarrow D) \geq \text{conf}(AB \rightarrow CD) \geq \text{conf}(A \rightarrow BCD)$$

$$\sigma(ABCD)$$

$$\sigma(ABCD)$$

$$\sigma(ABCD)$$

$$\sigma(ABC)$$

$$\sigma(AB)$$

$$\sigma(A)$$

- Given $X \rightarrow Y$ or $\{X, Y\}$, information needed to compute interestingness can be obtained from a contingency table

Contingency table

	Y	\bar{Y}	
X	f_{11}	f_{10}	f_{1+}
\bar{X}	f_{01}	f_{00}	f_{0+}
	f_{+1}	f_{+0}	N

f_{11} : support of X and Y
 f_{10} : support of X and \bar{Y}
 f_{01} : support of X and \bar{Y}
 f_{00} : support of \bar{X} and \bar{Y}

Used to define various measures
 support, confidence, Gini, entropy, etc.

- X and Y are independent:
 - $\text{Confidence}(X \rightarrow Y) = \text{support}(Y)$
 - $P(Y|X) = P(Y)$
 - $P(X,Y) = P(X) \times P(Y)$
- X & Y are positively correlated: $P(X,Y) > P(X) \times P(Y)$
- X & Y are negatively correlated: $P(X,Y) < P(X) \times P(Y)$
- Those can be measured by Lift or Interest:

$$\begin{aligned}
 \text{Lift} &= \frac{P(Y|X)}{P(Y)} \\
 \text{Interest} &= \frac{P(X,Y)}{P(X)P(Y)}
 \end{aligned}
 \quad \left. \begin{array}{l} \text{lift is used for rules while} \\ \text{interest is used for itemsets} \end{array} \right\}$$

- Simpson's Paradox: observed relationship in data may be influenced by the presence of other confounding factors (hidden variables). Hidden variables may cause the observed relationship to disappear or reverse its direction, so proper stratification is needed to avoid generating spurious patterns, e.g.

- Question

Hospital	Recovered	Recovered	
A	99	81	180
B	54	66	120
	153	147	300

$$\text{Conf}(\{\text{Hospital}=A\} \rightarrow \{\text{Recovered}\}) = 99/180 = 55\%$$

$$\text{Conf}(\{\text{Hospital}=B\} \rightarrow \{\text{Recovered}\}) = 54/120 = 45\%$$

=> Patient who stay at hospital A are more likely to be recovered?

- Answer:

Patient	Hospital	Recovered	Recovered	Total
Young	A	1	9	10
	B	4	30	34
Old	A	98	72	170
	B	50	36	86

Young Patients

$$\text{Conf}(\{\text{Hospital}=A\} \rightarrow \{\text{Recovered}\}) = 1/10 = 10\%$$

$$\text{Conf}(\{\text{Hospital}=B\} \rightarrow \{\text{Recovered}\}) = 4/34 = 11.8\%$$

Old Patients

$$\text{Conf}(\{\text{Hospital}=A\} \rightarrow \{\text{Recovered}\}) = 98/170 = 57.7\%$$

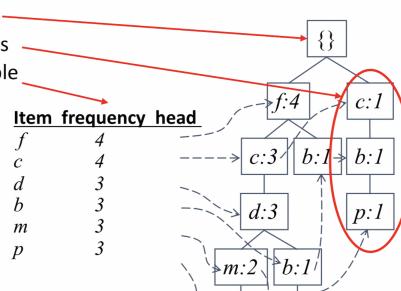
$$\text{Conf}(\{\text{Hospital}=B\} \rightarrow \{\text{Recovered}\}) = 50/86 = 58.1\%$$

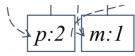
FP-Tree and FP-Growth for Data Structure

- FP-Tree

- Components

- One root: labeled as "null"
- A set of item prefix subtrees
- A frequent-item header table





- Construction

- Sort items by the order in the header table

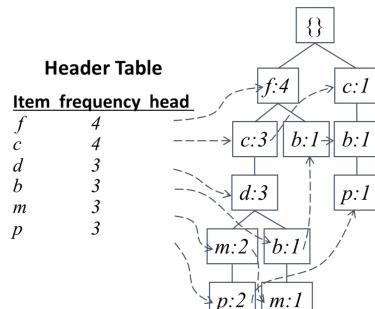
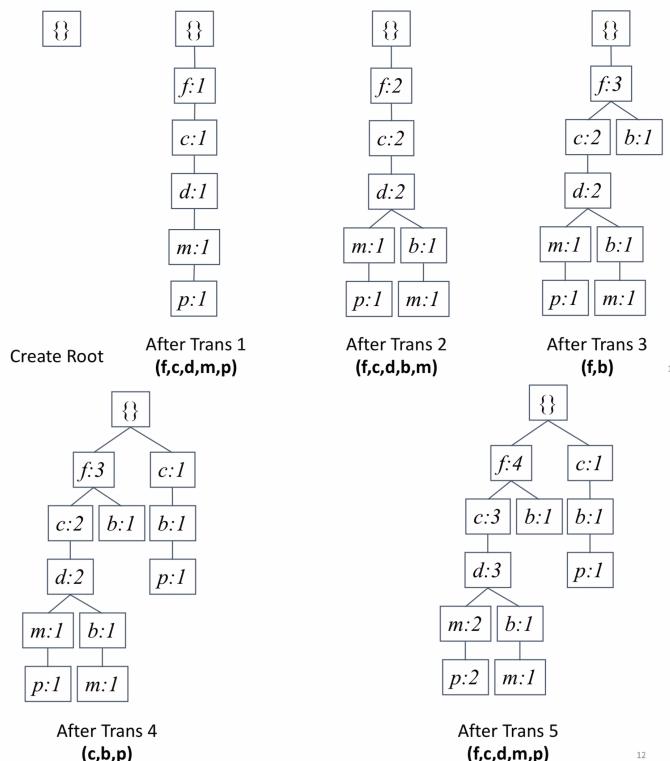
TID	Items Bought	(Ordered) Frequent Items
1	f,d,c,a,g,i,m,p	f,c,d,m,p
2	d,b,c,f,l,m,o	f,c,d,b,m
3	b,f,h,j,o	f,b
4	b,c,k,s,p	c,b,p
5	d,f,c,e,l,p,m,n	f,c,d,m,p

minimum support threshold (count) is set to 3

- First Scan count and sort

$$L = \{(f:4), (c:4), (d:3), (b:3), (m:3), (p:3)\}$$

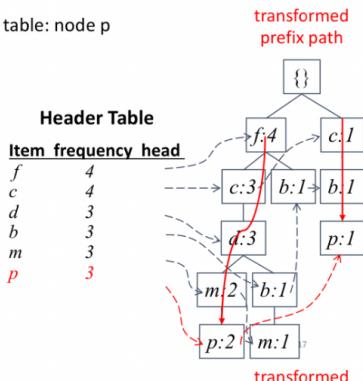
- count the frequencies of each item
- collect length 1 frequent items, then sort them in support/count descending order into L , frequent item list
- Second Scan create the tree and header table



- Create the root, label it as "null"
 - For each transaction
 - select and sort the frequent items in Trans
 - increase nodes count or create new node
 - If prefix nodes already exist, increase their counts by 1
 - If no prefix nodes, create it and set count to 1.
 - build the item header table: nodes with the same item name are linked in sequence via node links
- The size of the FP tree depends on how the items are ordered
 - Best case scenario : all transactions contain the same set of 1 path in the FP tree

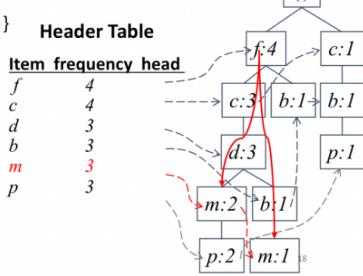
- Worst case scenario : every transaction has a unique set of items (no items in common)
- FP-Growth: recursively grow frequent patterns by pattern and database partition

- Start from the bottom of the header table: node p
- Two paths
- p's conditional pattern base
 - { (f:2, c:2, d:2, m:2), (c:1, b:1) }
- p's conditional FP tree
 - Only one branch (c:3) \rightarrow pattern (cp: 3)
- Patterns: (p:3), (cp:3)



- Continue with node m
- Two paths

- m's conditional pattern base
 - { (f:2, c:2, d:2), (f:1, c:1, d:1, b:1) }



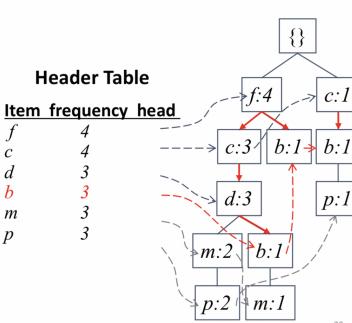
- Call mine(< f:3, c:3, d:3 > | m)
- Patterns:
 - (m:3)
 - See next slide

- node d:
 - (dm:3)
 - call mine(<f:3, c:3>|dm)
 - (cdm:3)
 - call(<f:3>|cdm) \rightarrow (fcdm:3), (fdm:3)
- node c:
 - (cm:3)
 - call mine(<f:3>|cm)
 - (fcm:3)
- node f: (fm:3)
- All the patterns: (m:3, dm:3, cm:3, fm:3, cdm:3, fdm:3, fcm:3, fcdm:3)
- A single path FP-Tree can be mined by outputting all the combination of the items in the path.

- Continue with node b

- Three paths

- b 's conditional pattern base
 - { (f:1, c:1, d:1), (f:1), (c:1) }
- b 's conditional FP tree
 - \emptyset
- Patterns: (b:3)



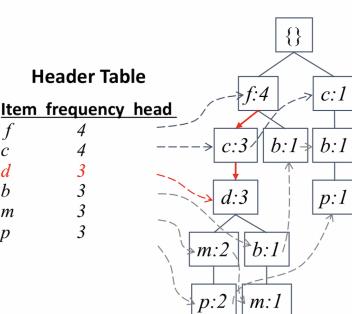
20

- Continue with node d

- One path

- d 's conditional pattern base
 - { (f:3, c:3) }
- d 's conditional FP tree
 - { (f:3, c:3) }

- Patterns:
 - (d:3)
 - (cd:3)
 - (fd:3)
 - (fcd:3)

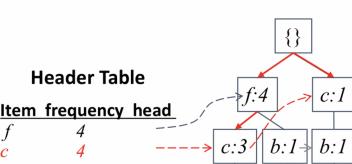


21

- Continue with node c

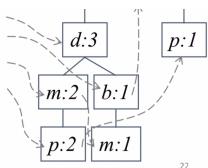
- Two paths

- c 's conditional pattern base
 - { (f:3) }
- c 's conditional FP tree
 - { (f:3) }



- Patterns:
 - (c:4)
 - (fc:3)

d 3
 b 3
 m 3
 p 3



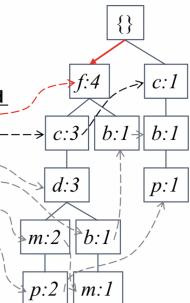
22

- Continue with node f
- One path
- f's conditional pattern base
 - \emptyset
- f's conditional FP tree
 - \emptyset
- Patterns:
 - (f:4)

Header Table

Item frequency head

f 4
 c 4
 d 3
 b 3
 m 3
 p 3



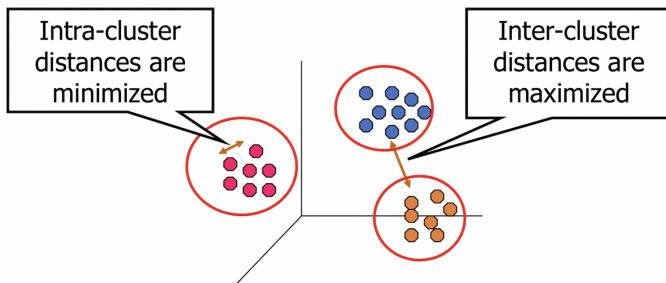
23

Item	Conditional Pattern Base	Conditional FP-Tree
p	{(f:2, c:2, d:2, m:2), (c:1, b:1)}	{(c:3)} p
m	{(f:2, c:2, d:2), (f:1, c:1, d:1, b:1)}	{(f:3, c:3, d: 3)} m
b	{(f:1, c:1, d:1), (f:1), (c:1)}	\emptyset
d	{(f:3, c:3)}	{(f:3, c:3)} d
c	{(f:3)}	{(f:3)} c
f	\emptyset	\emptyset

1. For each frequent item, construct its conditional pattern-base, and then its conditional FP-tree
2. Repeat the process on each newly created conditional FP-tree
3. Until the resulting FP-tree is empty, or it contains only one path (single path will generate all the combinations of its sub-paths, each of which is a frequent pattern)

Cluster Analysis

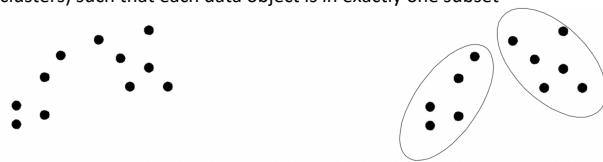
is finding groups of objects such that



1. the objects in a group will be similar (or related) to one another
2. different from (or unrelated to) the objects in other groups

Types of Clustering (i.e. a set of clusters)

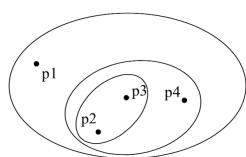
- Partitional versus Hierarchical
 - Partitional Clustering: a division of data objects into non overlapping subsets (clusters) such that each data object is in exactly one subset



Original Points

A Partitional Clustering

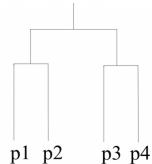
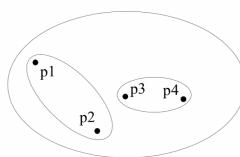
- Hierarchical clustering: a set of nested clusters organized as a hierarchical tree



Traditional Hierarchical Clustering

Traditional Dendrogram

|

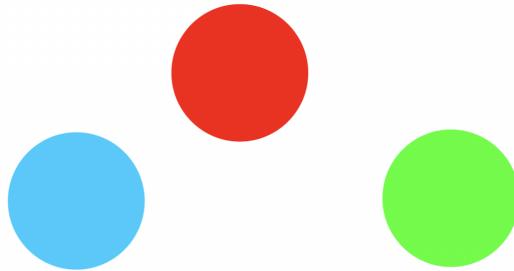


Non-traditional Hierarchical Clustering Non-traditional Dendrogram

- Agglomerative: start with the points as individual clusters, and at each step, merge the closest pair of clusters until only one cluster (or k clusters) is left
- Divisive: start with one, all-inclusive cluster, and at each step, split a cluster until each cluster contains an individual point (or there are k clusters)
- Exclusive versus non-exclusive: in non-exclusive clustering, points may belong to multiple clusters/multiple classes or could be 'border' points
- Fuzzy (one type of non-exclusive) versus non-fuzzy clustering: in fuzzy clustering, a point is with some weight between 0 and 1 indicating strength of belief in the event in question and weights must sum to 1 for that event, e.g. the point belongs to cluster a
- Partial versus complete: in partial clustering, we only want to cluster some of the data

Types of Cluster

- Well-separated: a cluster is a set of points such that any point in a cluster is closer (or more similar) to every other point in the cluster than to any point not in the cluster



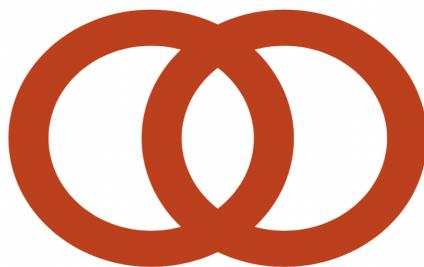
3 well-separated clusters

- Centre-based: a cluster is a set of objects such that an object in a cluster is closer (more similar) to the prototype or "centre" of a cluster, than to the centre of any other cluster, where the centre of a cluster is often
 - a centroid, the average of all the points in the cluster
 - a medoid, the most "representative" point of a cluster



4 center-based clusters

- Shared Property or Conceptual: clusters that share some common property or represent a particular concept



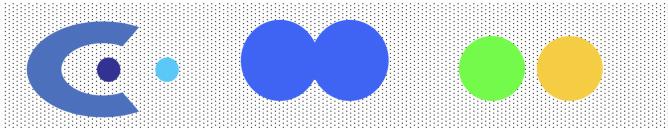
2 Overlapping Circles

- Contiguity-based (nearest neighbour or transitive): a cluster is a set of points such that a point in a cluster is closer (or more similar) to one or more other points in the cluster than to any point not in the cluster.



8 contiguous clusters

- Density-based: a cluster is a dense region of points, which is separated by low-density regions, from other regions of high density. It's used when the clusters are irregular or intertwined, and when noise and outliers are present.



6 density-based clusters

- Described by an Objective Function
 - Global vs. Local
 - Hierarchical clustering algorithms typically have local objectives
 - Partitional algorithms typically have global objectives
 - Parameterized vs. Mixture
 - Parameters for the model are determined from the data.
 - Mixture models assume that the data is a 'mixture' of a number of statistical distributions.
 - Minimize vs. Maximize in the context of a weighted graph
where the nodes are the points being clustered, and the weighted edges represent the proximities between points
 - minimize the edge weight between clusters and maximize the edge weight within clusters [if your objective function uses 1 to represent an exact match/similarity and 0 to represent no similarity]
 - maximize the edge weight between clusters and minimize the edge weight within clusters [if your objective function uses 0 to represent an exact match/similarity]
- Prototype-based clusters

Characteristics of Input Data

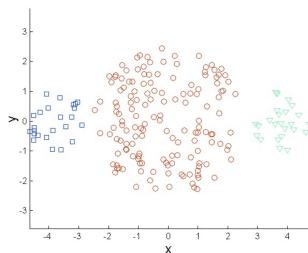
- Type of proximity or density measure
- Sparseness dictates type of similarity and efficiency
- Attribute type dictates type of similarity
- Type of Data dictates type of similarity and other characteristics, e.g., autocorrelation
- Noise and Outliers: often interfere with the operation of the clustering algorithm
- Clusters of differing sizes, densities, and shapes

K-means Clustering

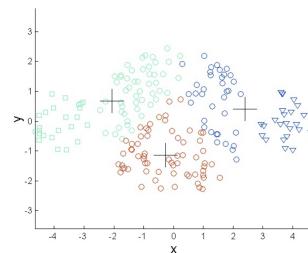
- Pre-processing
 - Normalize the data
 - Eliminate outliers

- 1: Select K points as the initial centroids.
- 2: **repeat**
- 3: Form K clusters by assigning all points to the closest centroid.
- 4: Recompute the centroid of each cluster.
- 5: **until** The centroids don't change

- 1. Initial centroids are often chosen randomly, so clusters produced vary from one to another
- 2. Time Complexity: $O(l * n * K * d)$ where n = number of points, K = number of clusters, l = number of iterations, d = number of attributes
- 3. 'Closeness' is measured by Euclidean distance, cosine similarity, correlation, SSE etc.
- 4. The centroid is (typically) the mean of the points in the cluster
- 5. Most of the convergence happens in the first few iterations for the common similarity measures above, and often the stopping condition is changed to 'Until relatively few points change clusters'
- Post-processing
 - Eliminate empty clusters and small clusters that may represent outliers
 - Split 'loose' clusters, i.e., clusters with relatively high SSE
 - Merge clusters that are 'close' and that have relatively low SSE
 - These steps can be used multiple times during the clustering process
- K-means has problems when clusters are of differing
 - Sizes

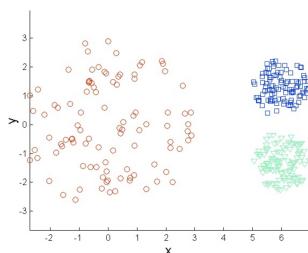


Original Points

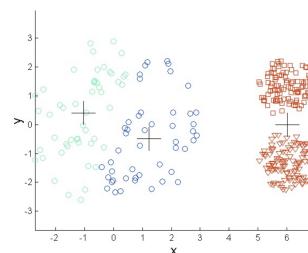


K-means (3 Clusters)

- Densities

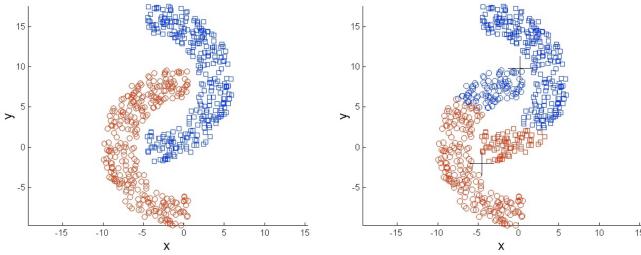


Original Points



K-means (3 Clusters)

- Non-globular shapes



Original Points

K-means (2 Clusters)

- One solution is to use many clusters, find parts of clusters, and put them together

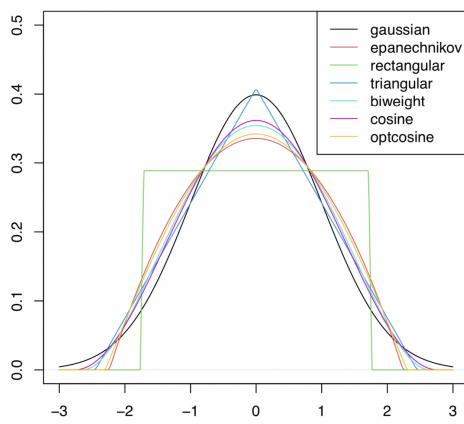
Mixture-based Clustering

is to fit a linear combination of probabilistic distributions, e.g. normal distributions,

$$f(x) = \pi_1 f_1(x) + \cdots + \pi_K f_K(x)$$

to the data, where each component represents a cluster, associated with h a mixing proportion π_j , subject to $\pi_j > 0$ and $\sum_{j=1}^K \pi_j = 1$.

- Similar to density estimation, a kernel function $K(x)$ (usually one that is symmetric about 0 and has a unit variance) is used to provide weights to observations in the neighbourhood of a point and a smoothing parameter h , known as the bandwidth, to control the size of the neighbourhood,



- Let's denote

$$K_h(x) = \frac{1}{h} K\left(\frac{x}{h}\right)$$

if $K(x)$ is the standard normal distribution, then $K_h(x)$ is the normal distribution with mean 0 and variance h^2

- The true distribution is denoted by f and transformed from the equation above
- For normal mixtures:
 - Equal variances (homoscedastic): All normal components share an identical variance σ^2 (but may have different means μ_j).
 - Varying variances (heteroscedastic): Each normal component has its own variance σ_j^2 (and mean μ_j).
- An observation x is of cluster j , if $\pi_j f_j(x)$ is the largest, i.e., the same rule we used for classification.
- Invariant of data scaling

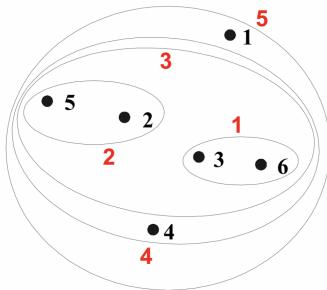
Agglomerative Clustering

- Compute the proximity matrix
 - Let each data point be a cluster
 - Repeat**
 - Merge** the two closest clusters
 - Update** the proximity matrix
 - Until** only a single cluster remains
- Complexity
 - $O(N^2)$ space since it uses the proximity matrix. • N is the number of points.
 - $O(N^3)$ time in many cases
 - There are N steps and at each step the size, N^2 , proximity matrix must be updated and searched
 - Complexity can be reduced to $O(N^2 \log(N))$ time with some approaches
 - The key operation is the computation of the proximity of two clusters, i.e. this depends how we prefer the clusters to be linked together, hence the notion of linkage, e.g.

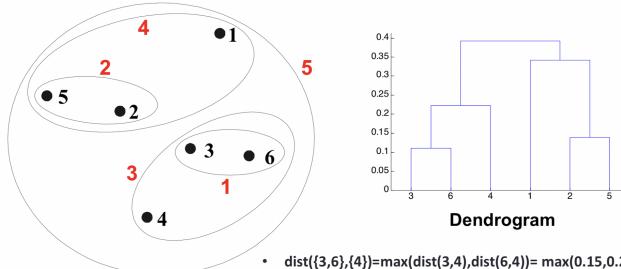
	p1	p2	p3	p4	p5	p6
p1	0.00	0.24	0.22	0.37	0.34	0.23
p2	0.24	0.00	0.15	0.20	0.14	0.25
p3	0.22	0.15	0.00	0.15	0.28	0.11

p4	0.37	0.20	0.15	0.00	0.29	0.22
p5	0.34	0.14	0.28	0.29	0.00	0.39
p6	0.23	0.25	0.11	0.22	0.39	0.00

- Single Linkage/Min: the proximity of two clusters is based on the two closest points in the different clusters, hence determined by one pair of points, i.e., by one link in the proximity graph



- Pro: can handle non-elliptical shapes
- Con: sensitive to noise and outliers
- Complete Linkage/Max: the proximity of two clusters is based on the two most distant points in the different clusters, hence determined by all pairs of points in the two clusters



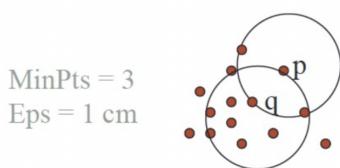
- dist({3,6},{4})=max(dist(3,4),dist(6,4))= max(0.15,0.22)=0.22.
- dist({3,6},{2,5})=max(dist(3,2),dist(6,2),dist(3,5),dist(6,5))= max(0.15,0.25,0.28,0.39)=0.39.
- dist({3,6},{1})=max(dist(3,1),dist(6,1))= max(0.22,0.23)=0.23.
- Pro: less susceptible to noise and outliers
- Cons: tends to break large clusters; biased towards globular clusters
- Average Linkage/Group Average ~ Ward's Method with the squared distance: the proximity of two clusters is the average of pairwise proximity between points in the two clusters:

$$\text{proximity}(\text{Cluster}_i, \text{Cluster}_j) = \frac{\sum_{\substack{p_i \in \text{Cluster}_i \\ p_j \in \text{Cluster}_j}} \text{proximity}(p_i, p_j)}{|\text{Cluster}_i| \times |\text{Cluster}_j|}$$

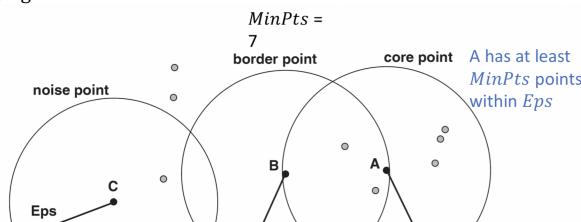
- Pros: compromise between Single and Complete Link Strengths; less susceptible to noise and outliers
- Cons: biased towards globular clusters

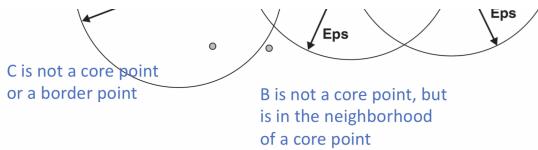
Density-Based Spatial Clustering of Applications with Noise (DBSCAN) Clustering

- Select a point p .
- Retrieve all points density reachable from p wrt. Eps and $MinPts$.
- If p is a core point, a cluster is formed.
- If p is a border point, no points are density reachable from p and DBSCAN visits the next point of the database.
- Continue the process until all the points have been processed.
- Parameters
 - Eps : Maximum radius of the neighbourhood
 - $MinPts$: Minimum number of points in an Eps neighbourhood
 - $NE(p)$: $\{q | dist(p, q) \leq Eps\}$ where the point q directly density reachable from p iff $q \in NE(p)$ and p is a core object p , e.g.

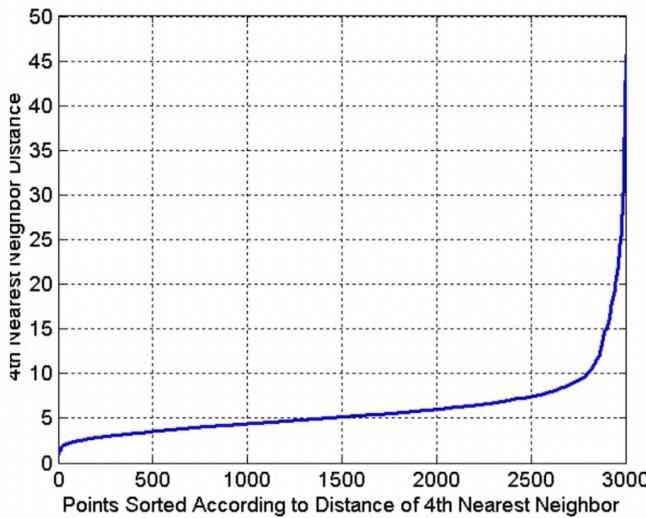


- A point is, e.g.





- a core point if it has at least a specified number of points (*MinPts*) within *Eps*
 - these are points that are at the interior of a cluster
 - counts the point itself
- a border point is not a core point, but is in the neighbourhood of a core point
- a noise point is any point that is not a core point or a border point
- Determining *EPS* and *MinPts* - for points in a cluster, their *k*th nearest neighbours are at close distance, and for noise points, they have the *k*th nearest neighbour at farther distance. So, plot sorted distance of every point to its *k*th nearest neighbour, e.g.

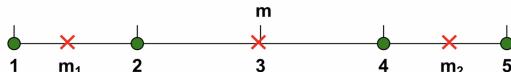


- Pros & Cons
 - Pros: is resistant to Noise and can handle clusters of different shapes and sizes
 - Cons: sensitive to varying densities and high-dimensional data

Cluster Validity

- External Index: Used to measure the extent to which cluster labels match externally supplied class labels, e.g. Entropy
- Relative Index: Used to compare two different clustering methods, e.g. SSE
- Internal Index: Used to measure the goodness of a clustering

$$WSS + BSS = \text{constant} \quad WSS = \sum_i \sum_{x \in C_i} (x - m_i)^2 \quad BSS = \sum_i |C_i|(m - m_i)^2$$



K=1 cluster: $WSS = (1 - 3)^2 + (2 - 3)^2 + (4 - 3)^2 + (5 - 3)^2 = 10$
 $BSS = 4 \times (3 - 3)^2 = 0$
 $Total = 10 + 0 = 10$

K=2 clusters: $WSE = (1 - 1.5)^2 + (2 - 1.5)^2 + (4 - 4.5)^2 + (5 - 4.5)^2 = 1$
 $BSS = 2 \times (3 - 1.5)^2 + 2 \times (4.5 - 3)^2 = 9$
 $Total = 1 + 9 = 10$

- Cluster cohesion measures how closely related are objects in a cluster, which is measured by the within cluster sum of squares (WSS), i.e. SSE -> the lower, the better

$$WSS = \sum_i \sum_{x \in C_i} (x - m_i)^2$$

- $|C_i|$ is the size of cluster *i*

- Cluster separation measures how distinct or well-separated a cluster is from other clusters, which is measured by the sum of the weights between nodes in the cluster and nodes outside the cluster (BSS) -> the higher, the better

$$BSS = \sum_i |C_i|(m - m_i)^2$$

- $|C_i|$ is the size of cluster *i*

Anomaly/Outlier Detection

Noise vs. Outliers vs. Anomaly

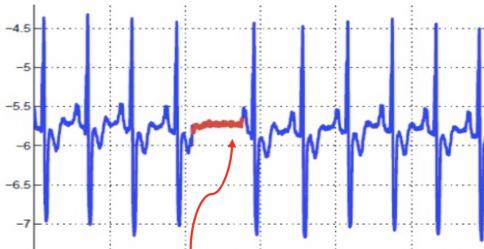
- Noise is random error or variance in a measured variable and should be removed before outlier detection
- Outlier: a data object that deviates significantly from the normal objects as if it were generated by a different mechanism, e.g. unusual credit card purchase, sports: Michael

Jordon, Wayne Gretzky, etc. Through treatment of outliers, it leads to improving the model accuracy.

- Anomaly is the outlier that can be explained by few features (may be new features). Through Anomaly Detection, understanding the pattern of anomalies, may lead to new findings (a new different model) or also, lead to new features that can be introduced in the existing model.

Types of Outliers

- Object O_g is a global outlier (or point anomaly) if it significantly deviates from the rest of the data set, e.g. intrusion detection in computer networks
- Object O_c is a contextual outlier (or conditional outlier) if it deviates significantly based on a selected context, can be also viewed as a generalization of local outliers—whose density significantly deviates from its local area
 - Contextual attributes: define the context, e.g. time & location
 - Behavioural attributes: characteristics of the object, used in outlier evaluation, e.g., temperature
- A subset of data objects are collective outliers if they collectively deviate significantly from the whole data set, even if the individual data objects may not be outliers, e.g. intrusion detection, when several computers keep sending denial-of-service packages to each other, or heartbeat monitoring



Based on assumptions about normal data and outliers

- Statistical methods (also known as model-based methods) assume that the normal data follow some statistical model (a stochastic model), and the data not following the model are outliers.
 - Assumption: data objects are generated by a stochastic process (a generative model)
 - Idea: learn a generative model fitting the given data set, and then identify the objects in low probability regions of the model as outliers
 - Approaches
 - Parametric method: the probability density function of the parametric distribution $f(x|\theta)$ gives the probability that object x is generated by the distribution: the smaller this value, the more likely x is an outlier, e.g.
 - Detecting Univariate Outliers Based on Normal Distribution

Avg. temp.: {24.0, 28.9, 28.9, 29.0, 29.1, 29.1, 29.2, 29.2, 29.3, 29.4}
Use the maximum likelihood method to estimate μ and σ^2

$$\ln \mathcal{L}(\mu, \sigma^2) = \sum_{i=1}^n \ln f(x_i | (\mu, \sigma^2)) = -\frac{n}{2} \ln(2\pi) - \frac{n}{2} \ln \sigma^2 - \frac{1}{2\sigma^2} \sum_{i=1}^n (x_i - \mu)^2$$

Taking derivatives with respect to μ and σ^2 , we derive the following maximum likelihood estimates

$$\hat{\mu} = \bar{x} = \frac{1}{n} \sum_{i=1}^n x_i \quad \hat{\sigma}^2 = \frac{1}{n} \sum_{i=1}^n (x_i - \bar{x})^2$$

For the above data with $n = 10$, we have $\hat{\mu} = 28.61$ and $\hat{\sigma} = \sqrt{2.29} = 1.51$
Then $(24 - 28.61) / 1.51 = -3.04 < -3$, **24 is an outlier** since

$\mu \pm 3\sigma$ region contains 99.7% data

- Detection of Multivariate Outliers: transform the multivariate outlier detection task into a univariate outlier detection problem
 1. Calculate the mean vector from the multivariate data
 2. For each object o , calculate the Mahalanobis distance from o to \bar{o} : $MDist(o, \bar{o}) = (o - \bar{o})^T S^{-1}(o - \bar{o})$ where S is the covariance matrix
 3. Detect outliers in the transformed univariate data set, $\{MDist(o, \bar{o}) | o \in D\}$ using a statistical test, e.g., Grubb's test. If $MDist(o, \bar{o})$ is determined to be an outlier, then o is regarded as an outlier
- Non-parametric: not assume an a-priori statistical model and determine the model from the input data, event not completely parameter free but consider the number and nature of the parameters are flexible and not fixed in advance, e.g. histogram and kernel density estimation
- Proximity based methods
 - Intuition: Objects that are far away from the others are outliers
 - Assumption: The proximity of an outlier deviates significantly from that of most of the others in the data set
 - Approaches
 - Distance-based outlier detection: An object o is an outlier if its neighbourhood does not have enough other points, e.g.

$$\frac{\| \{o' | dist(o, o') \leq r\} \|}{\|D\|} \leq \pi$$

- r is a user-specified distance threshold
- o' refers to the other objects within the distance threshold
- D refers to the neighbourhood

- π is a fraction threshold
- Density-based outlier detection: An object o is an outlier if its density is relatively much lower than that of its neighbours, e.g.
 - o is an outlier if k -distance neighbourhood of o $N_k(o) = \{o' \mid o' \text{ in } D, \text{dist}(o, o') \leq \text{dist}_k(o)\}$ is significantly smaller where k -distance of an object o , $\text{dist}_k(o)$: distance between o and its k -th nearest neighbour (NN) -> global outlier
 - Local Outlier Factor (LOF) -> local outlier, e.g. consider the following 4 data points: $a(0, 0)$, $b(0, 1)$, $c(1, 1)$, $d(3, 0)$, set $k=2$ and use Manhattan Distance

$$\text{dist}(a, b) = 1$$

$$\text{dist}(a, c) = 2$$

$$\text{dist}(a, d) = 3$$

1.

$$\text{dist}(b, c) = 1$$

$$\text{dist}(b, d) = 3+1=4$$

$$\text{dist}(c, d) = 2+1=3$$

$\text{dist}_k(o)$: distance between o and its k -th NN

- 2.
- $\text{dist}_2(a) = \text{dist}(a, c) = 2$ (c is the 2nd NN)
 - $\text{dist}_2(b) = \text{dist}(b, a) = 1$ (a/c is the 2nd NN)
 - $\text{dist}_2(c) = \text{dist}(c, a) = 2$ (a is the 2nd NN)
 - $\text{dist}_2(d) = \text{dist}(d, a) = 3$ (a/c is the 2nd NN)

k-distance neighbourhood of o ,

$N_k(o) = \{o' \mid o' \text{ in } D, \text{dist}(o, o') \leq \text{dist}_k(o)\}$

3. $N_2(a) = \{b, c\}$

$N_2(b) = \{a, c\}$

$N_2(c) = \{b, a\}$

$N_2(d) = \{a, c\}$

- Reachability distance from o' to o :

$$\text{reachdist}_k(o \leftarrow o') = \max\{\text{dist}_k(o), \text{dist}(o, o')\}$$

- Local reachability density of o :

$$\text{lrd}_k(o) = \frac{\|N_k(o)\|}{\sum_{o' \in N_k(o)} \text{reachdist}_k(o' \leftarrow o)}$$

4. • Example, object a, b, c, d

$$\text{lrd}_2(a) = \|N_2(a)\| / (\text{reachdist}_2(b \leftarrow a) + \text{reachdist}_2(c \leftarrow a)) = 2/(1+2) = 0.667$$

$$\text{lrd}_2(b) = \|N_2(b)\| / (\text{reachdist}_2(a \leftarrow b) + \text{reachdist}_2(c \leftarrow b)) = 2/(2+2) = 0.5$$

$$\text{lrd}_2(c) = \|N_2(c)\| / (\text{reachdist}_2(b \leftarrow c) + \text{reachdist}_2(a \leftarrow c)) = 2/(1+2) = 0.667$$

$$\text{lrd}_2(d) = \|N_2(d)\| / (\text{reachdist}_2(a \leftarrow d) + \text{reachdist}_2(c \leftarrow d)) = 2/(3+3) = 0.33$$

- LOF (Local outlier factor) of an object o : the average of the ratio of local reachability of o and those of o 's k nearest neighbors

$$\text{LOF}_k(o) = \frac{\sum_{o' \in N_k(o)} \frac{\text{lrd}_k(o')}{\text{lrd}_k(o)}}{\|N_k(o)\|} = \sum_{o' \in N_k(o)} \text{lrd}_k(o') \cdot \sum_{o' \in N_k(o)} \text{reachdist}_k(o' \leftarrow o)$$

5. • The lower the local reachability density of o , and the higher the local reachability density of the kNN of o , the higher LOF

- This captures a local outlier whose local density is relatively low comparing to the local densities of its kNN

$$\text{LOF}_2(a) = (\text{lrd}_2(b) + \text{lrd}_2(c)) * (\text{reachdist}_2(b \leftarrow a) + \text{reachdist}_2(c \leftarrow a)) = (0.5+0.667) * (1+2) = 3.5$$

$$\text{LOF}_2(b) = (\text{lrd}_2(a) + \text{lrd}_2(c)) * (\text{reachdist}_2(a \leftarrow b) + \text{reachdist}_2(c \leftarrow b)) = (0.667+0.667) * (2+2) = 5.3$$

$$\text{LOF}_2(c) = (\text{lrd}_2(b) + \text{lrd}_2(a)) * (\text{reachdist}_2(b \leftarrow c) + \text{reachdist}_2(a \leftarrow c)) = (0.5+0.667) * (1+2) = 3.5$$

$$\text{LOF}_2(\text{d}) = (\text{lrd}_2(a) + \text{lrd}_2(c)) * (\text{reachdist}_2(a \leftarrow d) + \text{reachdist}_2(c \leftarrow d)) = (0.667+0.667) * (3+3) = 8.0$$

- Clustering based methods

- Cases

- Not belong to any cluster -> a density-based clustering method such as DBSCAN to identify an outlier, e.g. identify animals not part of a flock
- Far from its closest cluster -> a partition clustering method such as k-means, e.g. intrusion detection
- Belongs to a small or sparse cluster -> find outliers based on cluster-based local outlier factor (CBLOF)

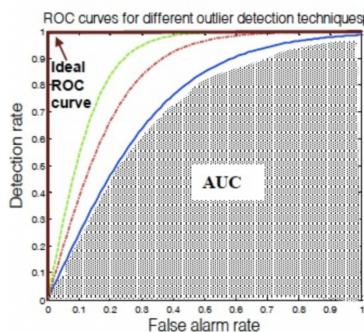
- Pros & Cons

- Pros
 - Detect outliers without requiring any labeled data
 - Work for many types of data
 - Clusters can be regarded as summaries of the data
 - Once the cluster are obtained, only need to compare any object against the clusters to determine whether it is an outlier (fast)
- Cons
 - Effectiveness depends highly on the clustering method used
 - High computational cost to find clusters first, and this can be reduced by fixed-width clustering
 - A point is assigned to a cluster if the centre of the cluster is within a pre-defined distance threshold from the point
 - If a point cannot be assigned to any existing cluster, a new cluster is created, and the distance threshold may be learned from the training data under certain conditions

Based on whether user labelled examples of outliers can be obtained

- Supervised, e.g. train a classification model that can distinguish “normal” data from outliers
- Semi-supervised, e.g. using a clustering-based approach, find a large cluster, C, and a small cluster, C1 and treat all objects in C as normal. Then, train a classification model that can distinguish “normal” data from outliers.
- Unsupervised, e.g. isolation forest: A collection of tree that recursively partition the data set,
 - In each iteration of the process, a random feature is selected, the data is then split based on a randomly chosen value.
 - The value is between the minimum and maximum of the chosen feature. (Follows the intuition of random forest.)
 - Repeat the process until the entire data set is partitioned to form an individual tree in the forest.
 - Outliers generally form much shorter paths from the root than normal data points
- Evaluation

Confusion Matrix		Predicted Class	
Actual Class	Normal	Normal	Outlier
	Outlier	TN	FP
		FN	TP



- Recall (TP) : ratio between the number of correctly detected anomalies and the total number of anomalies, i.e., detection rate
- False alarm (FP) : ratio between the number of data records from normal class that are misclassified as anomalies and the total number of data records from normal class
- ROC Curve : a trade-off between detection rate and false alarm rate
- Area under the ROC curve (AUC): computed using a trapezoid rule

Data Stream

An algorithmic abstraction to support real-time analytics

- Requirements
 - Process an instance at a time, and inspect it (at most) once
 - Use a limited amount of time to process each instance
 - Use a limited amount of memory
 - Provide answer anytime (prediction, clustering, patterns)
 - Adapt to temporal changes (Concept drift)
- Types
 - Approximation: data stream is large and fast, so accept approximate solutions in order to use less time and memory
 - Classifier:
 - Change Management: receive the label for an example seen in the past, and use it for adjusting the model, that is, for training
 - Prediction: receive an unlabelled example and make a prediction for it based on its current model.

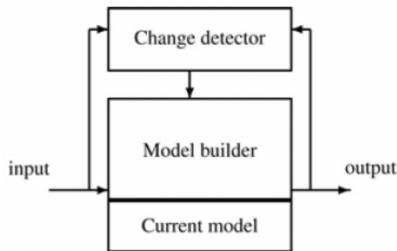
Approximation

- Sampling approach: for each item in the stream, process it with probability α , and ignore it with probability $1 - \alpha$, e.g. reservoir sampling for maintaining a fixed-size uniform random sample with the probability of k/t
 - Put the first k elements from the stream into the repository (Reservoir S)

- When the t-th element arrives where $t > k$, add it to reservoir S with probability $p = k/t$.
- Randomly remove an element from S if it is to be added
- Counting items and distinct items, e.g. the trivial solution keeps a standard integer counter and has the update operation just increment the counter
- Frequency problems: estimating the frequencies of all items is impossible to do in general in sublinear memory (some fixed, pre-set amount of memory), but often it suffices to have approximate counts for the most frequent items, e.g. lossy count: given a threshold ϵ , after reading any number of items from the stream, produce the set of heavy hitters where all items whose relative frequency exceeds ϵ
 1. Track items and counts.
 2. For each block of $1/\epsilon$ items, merge with stored items and counts.
 3. Decrement all counts by one, delete items with zero count.
- Counting within a sliding window to emphasize the most recent data and bound memory
- Merging sketches to count the number of 1s or 0s in the last k bit of a binary stream, e.g. exponential histogram for 1s
 - If the current bit is 0, then other changes are needed
 - If the current bit is 1, then
 1. Create a new bucket of size 1, for just this bit
 2. If there are now three buckets of size 1, combine the oldest two into a bucket of size 2
 3. If there are now three buckets of size 2, combine the oldest two into a bucket of size 4
 4. And so on and so forth ...

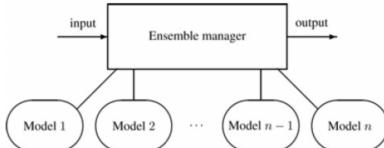
Change Management for Classifier

- Types of Change
 - Virtual drift occurs when $\Pr[X]$ changes but $\Pr[Y|X]$ remains unchanged.
 - Concept drift occurs when $\Pr[Y|X]$ changes, with or without changes in $\Pr[X]$.
- Strategies
 - Adaptive Estimators: monitoring a set of statistics from the stream and combining them into a model
 - Explicit Change Detectors for Model Revision



- One or more change detection algorithms run in parallel with the main model-building algorithm
- If significant change in the stream is detected, then activate a revision algorithm

- Model Ensemble



- Build complex classifiers out of simpler ones
- Single or several model-building algorithms are called at different times on different subsets of the data stream.
- An ensemble manager algorithm contains rules for creating, erasing, and revising the models in its ensemble.

- Types of Drift Detector

- Sequential analysis - Cumulative Sum (CUSUM): If $g_t > h$: declares change, reset $g_t, \mu, \sigma = 0$

$$g_0 = 0, g_t = \max(0, g_{t-1} + z_t - v)$$

- h is threshold
- v is drift speed
- $z_t = (x_t - \mu)/\sigma$, if μ and σ are not known a priori, they are estimated from the sequence itself
 - μ : the expected value of the x_t
 - σ : the standard deviation in "normal" conditions;

- Statistical Process Control - Drift Detection Method (DDM): Let P_t denote the error rate of the predictor at time t . Given the number of errors in a sample of t examples, its standard deviation at time t is given by:

$$S_t = \sqrt{P_t(1 - P_t)/t}$$

where the smallest value P_{min} of the error rates, observed up to time t , and the standard deviation S_{min} at that point are stored

- Warning is declared if $P_t + S_t \geq P_{min} + 2 \cdot S_{min}$: from this point on, new examples are stored in anticipation of a possible declaration of change.
- Change is declared if $P_t + S_t \geq P_{min} + 3 \cdot S_{min}$: the model induced by the learning method is discarded and a new model is built using the examples stored since the warning occurred. The values for P_{min} and S_{min} are reset.

- Monitoring two distributions - ADaptive sliding WINDOW (ADWIN) is a estimation algorithm based on exponential histograms where statistical test normally used is Hoeffding Bound and δ is error tolerance rate
 - No change: If W_0 and W_1 were generated from the same distribution (no change), then with probability at least $1-\delta$ the test says, "no change."
 - Change: If W_0 and W_1 were generated from two different distributions whose average differs by more than some quantity $\epsilon(W_0, W_1, \delta)$, then with probability at least $1-\delta$ the test says "change."
- Contextual approaches - ensemble of detectors as Diversity for Dealing with Drifts (DDD) acknowledges the diversity of different types of concept drift and tries to incorporate a variety of base classifiers in their ensemble.
- Evaluation
 - Mean time between false alarms (MTFA): Measures how often we get false alarms when there is no change. The false alarm rate (FAR) is defined as $1/\text{MTFA}$.
 - Mean time to detection (MTD): Measures the capacity of the learning system to detect and react to change when it occurs.
 - Missed detection rate (MDR): Measures the probability of not generating an alarm when there has been change.
 - Average run length ($\text{ARL}(\theta)$): Measures how long to wait before detecting a change after it occurs, which generalizes MTFA and MTD.

Prediction for Classifiers

- Evaluation Framework
 - Error estimation
 - Holdout: When data is so abundant that it is possible to have test sets periodically, then we can measure the performance on these holdout sets.
 - Interleaved test-then-train: Each individual example can be used to test the model before it is used for training, and from this, the accuracy can be incrementally updated.
 - Prequential: Like interleaved test-then-train but the idea that more recent examples are more important, using a sliding window or a decaying factor.
 - Interleaved chunks: Also, like interleaved test-then-train, but with chunks of data in sequence.
 - Evaluation performance measures, e.g. use a chance classifier to randomly assigns to each class the same number of examples as the classifier under consideration, then the Kappa statistic κ defined as follows:

$$\kappa = \frac{P_0 - P_c}{1 - P_c}$$

where P_0 : the classifier's prequential accuracy and P_c : the probability that a chance classifier makes a correct prediction (P_m : the probability of a majority class classifier makes a correct predication for the K_m statistic)

- $\kappa = 1$ if the classifier is always correct
- $\kappa = 0$ if its predictions coincide with the correct ones as often as those of a chance classifier

- Statistical significance validation
- A cost measure of the process
- Classifiers
 - The Majority Class algorithm is one of the simplest classifiers: it predicts the class of a new instance to be the most frequent class.
 - No-change classifier predicts the label for a new instance to be the true label of the previous instance. It exploits autocorrelation in the label assignments, which is very common.
 - Lazy Classifier consists of instances seen and predicts using the class label of the closest instances to the instance whose class label we want to predict.
 - Hoeffding Tree is a very fast decision tree (VFDT) algorithm for streaming data where instead of reusing instances, we wait for new instances to arrive

```

RECURSIVE_HUFFMAN_TREE
LET HT BE A TREE WITH SINGLE LEAF(ROOT);
INIT COUNTS nij AT ROOT;
for each example (x, y) in Stream do
  | HTGROW((x, y), HT);
end

def HTGROW((x, y), HT):
  | SORT (x, y) TO LEAF L USING HT ;
  | UPDATE COUNTS nij AT LEAF L ;
  | If examples seen so far at L are not all of the same class then
    |   | COMPUTE G FOR EACH ATTRIBUTE
    | end
  | If G(Best Attr.)- G(2nd best) >  $\sqrt{\frac{B^2 \ln(1/\delta)}{2n}}$  then → Use Hoeffding Bound to decide necessary
    |   | SPLIT LEAF ON BEST ATTRIBUTE
  | end
  | for each branch do
    |   | START NEW LEAF AND INITILIAZE COUNTS
  | end
end

```

- n : sample collected at a node
- r : range of the random variable, i.e., [max-min]
- δ : desired certainty, e.g., 99%
- Concept-adapting very fast decision tree (CVFDT) algorithm as an extension of VFDT to deal with concept drift, maintaining a model that is consistent with the instances stored in a sliding window.
- Hoeffding Adaptive Tree (HAT) replaces frequency counters by estimators, so it's parameter-free change detector + estimator with theoretical guarantees for subtree swap (ADWIN)
- Ensemble
 - Ozabag: it is difficult at first to draw a sample with replacement from the stream
The number of copies K of each of the n examples in the training set follows a binomial distribution and for large values of n , this binomial distribution tends to be a Poisson(1) distribution for bootstrap replicas
 - Adaptive Random Forest: it adds diversity to Ozabag by randomly selecting subsets of features for node splits

- It has one drift and warning detector per tree, which cause selective resets in response to drifts.
- It also allows training background trees, which start training if a warning is detected and replace the active tree if the warning escalates to a drift.
- Recurrent Concept Frameworks keep a library of classifiers that have been useful in the past but have been removed from the ensemble, for example because at some point they underperformed. Their accuracy is rechecked periodically in case they seem to be performing well again, and if so, they are added back to the ensemble.
- ADWIN Bagging for M models
- Leveraging Bagging
- Oza and Russell's Online Boosting