

DATA PREPROCESSING

ignore, manually, global central tendency, local central tendency, most probable value (reference)

1) Data Cleaning → missing values + noise handling + inconsistent data (eg. date format)

NOISE

- binning (local smoothing) → sort data and partition into bins, then smooth by bean → connecting a point neighbours
- outliers analysis → via clustering method
- filtering
- regression

- rectangular sliding average → replace the signal with the average of m adjacent points

• if the signal has a non-zero second-derivative ($f''(t) \neq 0$) → I will introduce a bias, but x -location of extremums are kept

- triangular smooth → correspond to two application of the rectangular filter (base 5 = 2 base 3) → more effective in reducing high frequency noise

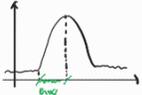
- Savitzky-Golay Filter → replace a point with a linear combination of adjacent $S_j = \sum_{n=-n_c}^{n_c} C_n y_{j+n}$
higher order polynomial

• preserve higher moments of the signal

• since least-squares involve matrix inversion (linear operation) → the coefficients C_n are themselves linear

smoothing ratio = $\frac{\text{window size}}{\text{coef. width}}$

• increasing the sr will increase SNR of the signal
 • $sr < 0.2$ if I want to measure height and width of a peak



→ smoothing for cosmetic reason or if the algorithm will be penalized (maximum location)

• but not for increasing accuracy, it will underestimate errors and smoothed noise may be mistaken for signal

$$f(x) = a_0 + a_1x + \dots + a_nx^n$$

$$A\vec{a} = \vec{y} \Rightarrow \vec{a} = (A^T A)^{-1} A^T \vec{y}$$

considering the window $\rightarrow A = \begin{bmatrix} 1 & x_1 & x_1^2 & \dots & x_1^n \\ \vdots & \vdots & \vdots & \dots & \vdots \\ 1 & x_m & x_m^2 & \dots & x_m^n \end{bmatrix}$

$$y_j = S_j = f(x_j) = a_0$$

$$\Rightarrow a_0 = \left[(A^T A)^{-1} A^T \vec{y} \right]_0$$

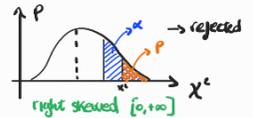
\vec{e}^T (first row)

2) Data Integration combining data from more sources → solve conflicts, redundancies (correlation analysis)

• **CHI-SQUARED TEST** → correlation between two categorical variables via the contingency table
 • non parametric • $H_0 \rightarrow$ two var are statistically independent

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

$= \frac{n_{11} \cdot n_{22} \cdot n_{33}}{n^3}$



α (significance level) → probability of rejecting H_0 when is actually true (concluding a dependency when there isn't)

• **PEARSON PRODUCT** (corr coefficient)

$$r_{AB} = \frac{\sum (a_i - \bar{a})(b_i - \bar{b})}{n \sigma_a \sigma_b} = \frac{\sum (a_i b_i) - n \bar{a} \bar{b}}{n \sigma_a \sigma_b}$$

• if $r_{A,B} > 0 \Rightarrow A, B$ are positively correlated (linear relationship)

• **CORRELATION**

$$a_i' = \frac{a_i - \bar{a}}{\sigma_a} \quad b_i' = \frac{b_i - \bar{b}}{\sigma_b} \Rightarrow \text{corr}(A, B) = \bar{a}' \cdot \bar{b}' = r_{AB} \cdot n$$

• **COVARIANCE**

$$\text{Cov}(X, Y) = E[(X - E[X])(Y - E[Y])] = E[XY] - E[X]E[Y]$$

$$= \frac{\sum (x_i - \bar{x})(y_i - \bar{y})}{n} = r_{A,B} \cdot \sigma_A \cdot \sigma_B$$

• X, Y independent $\Rightarrow \text{Cov}(X, Y) = 0$ - contrary only if $X, Y \sim N$

3) DATA REDUCTION find a new data representation with volume much smaller but comparable analytical results (reduce disk usage and computation time) also noise and irrelevant features

• **DIMENSIONALITY REDUCTION** → reduce the # of features of D
 → curse of dimensionality → data increasingly sparse (d less meaningful + # of data sample to cover the space grows exponentially)

→ PCA: finding a projection that captures the largest variation in data → **unsupervised**

normalize the dataset $A \cdot D \cdot \text{Dimen}$ → mean of dimension i

$$\text{Cov}(A, A) = \frac{AA^T}{n-1} \rightarrow \text{Bessel correction}$$

symmetric → diagonalizable + autoorthogonal; orthogonal; formula is valid since mean is null

Variance of projected data: $\frac{1}{n} \sum_{i=1}^n (\vec{u}^T \vec{x}_i)^2 = \vec{u}^T \Sigma \vec{u}$

$$\frac{\partial V}{\partial \vec{u}} = \frac{\partial}{\partial \vec{u}} \left[\vec{u}^T \Sigma \vec{u} + \lambda (1 - \vec{u}^T \vec{u}) \right] = 2 \Sigma \vec{u} - 2 \lambda \vec{u} = 0 \Rightarrow \Sigma \vec{u} = \lambda \vec{u} \Rightarrow \vec{u}^T \Sigma \vec{u} = \lambda$$

then compute the eigen-vectors of $\text{Cov}(A, A)$ and discard some with lower eigen-values → put them inside F
 $F = E^T A$
 • unsupervised + computationally heavy (matrix inversion)
 • near representation

→ **HEURISTIC SEARCH ATTRIBUTES**: 2^d possible attributes combinations → **greedy approach**

- best k attribute selection (given some metrics)
- best step wise attribute selection (select the first, then the second conditioned to the first)
- step wise attribute elimination (repeatedly remove the worst attributes)
- branch and bound (selection + elimination)

Can we Chi test with the attribute and the output or using the **normalized mutual information feature selection**

mutual info: $I(X, Y) = \sum_{x, y} p(x, y) \log \left(\frac{p(x, y)}{p(x)p(y)} \right) \rightarrow I(X, Y) = 0$ when independent

$$NI(f_i, f_j) = \frac{I(f_i, f_j)}{\min\{I(f_i), I(f_j)\}} \in [0, 1]$$

$$G = I(C, f_i) - \frac{1}{|S|} \sum_{f_j \in S} NI(f_i, f_j)$$

→ compensates to MI bias towards multivalued features

redundancy penalization

• supervised method
 • not just linear dependencies

→ **ATTRIBUTE CREATION** in order to get more informative/effective representation like with FFT

pseudoinverse $\min \|X\bar{w} - \bar{y}\|_2$
 $\frac{\partial f(\bar{w})}{\partial \bar{w}} = X^T(X\bar{w} - \bar{y}) = X^T X \bar{w} - X^T \bar{y} = 0$
 $\bar{w} = (X^T X)^{-1} X^T \bar{y}$

NUMEROSITY REDUCTION: reduce the # of data records or simplify their representation

PARAMETRIC METHODS → assume data fits some model, then store the parameters, discarding data

Regression analysis: least squared error to get the best fit of a dependent variable against one or more independent one

$w_i = \frac{\sum (x_i - \bar{x})(y_i - \bar{y})}{\sum (x_i - \bar{x})^2}$ $w_0 = \bar{y} - w_1 \bar{x}$ or with more variables: $\bar{w} = (X^T X)^{-1} X^T \bar{y}$
 • can get polynomial by changing X and adding the powers

NON PARAMETRIC METHODS → group in major families

- HISTOGRAM ANALYSIS → divide data in bucket and store avg (same width or same frequency/depth)
- CLUSTERING → store cluster representation (centroid and centroids)
- SAMPLING → choose a representative subset, simple random base poor performance if skew distribution → reduce the complexity of alg
 - w/w/o replacement or stratified sampling (partition dataset and sample accordingly) → eg. classes
- DATA CUBE AGGREGATION → each dimension store an attribute, cell of the cube store aggregated measures
 eg. date, location, product → store aggregated sales

DATA COMPRESSION → lossy vs. lossless

3) DATA TRANSFORMATION → map the entire set values of one attribute with a replacement set, where every point got mapped

smoothing, attribute construction, aggregation, normalization, discretization

min-max normalization → $v' = \frac{v - \min_A}{\max_A - \min_A} \cdot (\text{new_max}_A - \text{new_min}_A) + \text{new_min}_A$ Z-score normalization → $v' = \frac{v - \mu_A}{\sigma_A}$
 ↳ new range to be mapped

decimal scaling normalization → $v' = \frac{v}{10^J}$ J: smallest integer that $\max(|v'|) < 1$ $C = [-1, +1]$

DISCRETIZATION

attribute types: ^{unordered set} nominal, ^{ordered} ordinal, ^R numeric

can be merge (bottom-up) or split (top-down)

- BINNING (unsupervised, split) → equal width - equal depth (increases with the same # of samples)
- BY CLASSIFICATION (supervised, split) → using decision tree to find the optimal cut point → entropy of partitions
- BY CORRELATION (supervised, merge) → using a chi-merge algorithm → chi-squared between two neighboring intervals and the classes
 • if the class is independent, I should merge the intervals $df = \#classes - 1$
- BY CLUSTERING (unsupervised split/merge)

FAYYAD and IRANI approach optimal cut point → between two examples of different classes in the sorted values of A
 • I can take the midpoint

class entropy ← $H(S) = -\sum_{i=1}^K p(C_i, S) \cdot \log(p(C_i, S))$
 ↳ datasets ↳ output classes

- $H(S) = 0$ when I have only one class
- I want intervals with low class entropy → entropy as partitions quality measure

entropy of partitions ← $EP(A, T, S) = \frac{|S_1|}{N} H(S_1) + \frac{|S_2|}{N} H(S_2)$
 ↳ attribute ↳ cut point

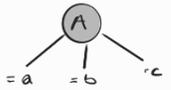
- Select the T_A with lowest EP → $S = S_1 \cup S_2$
- partition recursively S_1 and S_2 until a termination condition is met (no thresholds defined)

HIERARCHY GENERATION

→ from low level concept (eg. age) to higher level (eg. young, adult)
 • need to define a partial/total order (for numeric data I need to discretize first)

- street < city < state < country → schema level
- { Urbana, Chicago } < Illinois → explicit data grouping
- Automatic generation by occurrence counting → exception (day-week-month-year)

CLASSIFICATION → two step process: model construction + model usage / model testing (Eager Classifier)
 • supervised or unsupervised (went to understand 4 3 classes)



DECISION TREE → internal node does a test on an attribute, branches are the different outcomes and leaf hold a class values
 • no parameter setting, domain knowledge → perform automatically the feature selection process
 • termination when 1) no more sample 2) pure partition 3) no more attribute
 • very interpretable • induction complexity grows exponentially with the height
 • initially all attributes are at the cost

Greedy algorithm to find at each node the purest partition → different metrics

1) Information Gain $I(D) = - \sum p_i \cdot \log_2(p_i)$
 (Annotations: # bit to represent a class, class probability)

$I_A(D) = - \sum \frac{|D_j|}{|D|} I(D_j)$
 (Annotation: mean # bits to represent a class after splitting A)

$Gain(A) = I(D) - I_A(D)$
 (Annotation: reduction in info requirements by knowing the value of A)

- if can use this for **visualization** → sort A, then do a **binary** split for the midpoint with the lowest $I_A(D)$
 • infeasible for big datasets
- biased towards attribute with high # of values → eg. 10
 • the attribute with biggest gain will be selected

2) Gain Ratio $SplitInfo_A(D) = - \sum \frac{|D_j|}{|D|} \log_2 \frac{|D_j|}{|D|}$
 (Annotation: # of bits to express a partition)

→ if I have a lot of partition is bigger

$GainRatio_A(D) = \frac{Gain(A)}{SplitInfo(A)}$

select the highest

• new bias towards unbalanced split

3) Gini Index $gini(D) = 1 - \sum p_i^2$
 (Annotation: impurity's reduction)

• consider just **binary** splits
 $2^m - 2$

$gini_A(D) = \frac{|D_1|}{|D|} gini(D_1) + \frac{|D_2|}{|D|} gini(D_2)$
 (Annotation: twice # of classes in the splits)

$\Delta gini(A) = gini(D) - gini_A(D)$

• **Overfitting** when I have too many branches → poor generalization $A_{test} \ll A_{train}$ → overfitting

• **Generalization Error**

1) optimistic (ers = em)
 (Annotations: percentage, errors in D, |D|)

2) pessimistic $Ers = E_{TR} + \frac{0.5 \cdot N}{|D|}$ → # of leaf in decision tree (add check sample error for each leaf)

3) we test set to get the overall

4) we **prune** set to get generalization error of a leaf or a sub-tree

OCCAM RAZOR for models with same A I would prefer the simpler one, since the other probably captured also noise and address → more robust

• **Repetition** (attribute repeatedly tested, continuous) or **Redundation** (duplicate subtree) → can be reduced by **attribute construction** from sparsely represented trees

• **Pre-pruning** → early stopping, but I need some threshold

- if number of instances are less than k
- if class is independent of available features (chi-squared test)
- if expanding the node doesn't improve the impurity measures (info gain or gini index ...)

• **Post-pruning** → removing branches and substituting with a leaf of the most frequent class

REDUCED ERROR PRUNING (REP)

$gain = \# \text{ misclass}(T) - \# \text{ misclass}(v)$ → prune the largest gains until only negative ones remain
 (Annotations: on prune set, if pruned in v, recalculated on all non-leaf nodes)
 • bottom-up restriction: T can be pruned \iff doesn't contain sub-tree with lower error than v

COST COMPLEXITY PRUNING

$cost = \text{reconstruction error} + \beta \cdot \# \text{ leaves}$
 • start from the bottom, if the cost of the pruned is less \Rightarrow prune

- 1) bottom-up pruning if cost is less
- 2) order the trees at different steps by reconstruction error
- 3) Pick the tree with the smallest # of leaves within are std of error from the lowest cost tree

PESSIMISTIC PRUNING Like the above but with training error + some value derived from statistic

entropies calculations very efficient ←

RAIN FOREST → memory efficient algorithm • if training set doesn't fit in memory, with decision trees I need a lot of steps

AVC-set (attribute, value, class label) is the projection of the training dataset to the attribute x and class values. The **individual count** is aggregated
 • mention one at each node to describe n tuples at that node

AVC-group: Set of AVC-set for all the attributes at node n

AVC	Class	
	1	2
year	4	2
month	0	1
old	3	7

BOOST → use bootstrap sampling, where each fits into memory → get a **set of trees** that I use to get a unique T' that is very similar to what I would get normally
 (Annotation: uniform with replacement)

- Only two DS seen
- 1) small subset in order to get the underlying distribution
- 2) compute the final split with the info got from 1)
- I can use it for **incremental updates**

read partition by partition the DS

BAYESIAN CLASSIFICATION

$P(H|E) = \frac{P(E|H) \cdot P(H)}{P(E)}$

 (Posterior) = (Likelihood) * (Prior) / (Evidence)

 $P(E) = P(E|H) \cdot P(H) + P(E|\neg H) \cdot P(\neg H)$

 product of two gaussian is gaussian (conjugate prior)

prior can be estimated from data $P(C_i) = \frac{|C_i|}{|D|}$

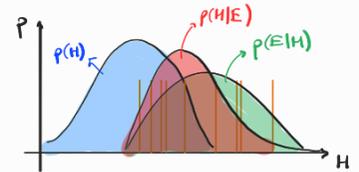
looking accuracy because dependencies exists

I want to find the $\max_{c_i} P(E|H) P(H)$

assuming independence of attributes

$P(X_k|C_j) = \prod_{k=1}^n P(x_k|C_j)$

- If X_k is continuous I will use a gaussian distribution $P(x_k|C_j) = g(x_k, \mu_j, \sigma_j) = \frac{1}{\sqrt{2\pi}\sigma} e^{-\frac{(x-\mu)^2}{2\sigma^2}}$
- I need to use the Laplacian Correction (adding 1 to each case)
- $P(H)$ reflects the likelihood, If we do not want to assume anything I can use non informative priors (uniform)

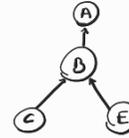


BAYESIAN BELF NETWORKS → represent attribute dependencies via a DAG. I can also add hidden variables to the model

CPT (conditional probability table) → that maps input probabilities to the one of the node

$P(x_1, \dots, x_n) = \prod_{i=1}^n P(x_i | \text{parents}(x_i))$

A node is selected as an output with the class label



$P(A, B, C, D, E) = P(A) \cdot P(B|A) \cdot P(C|B) \cdot P(D|B) \cdot P(E|D)$

	C_i, E	$C_i, \neg E$	$\neg C_i, E$	$\neg C_i, \neg E$
P_i	0.8	x	y	1-z
$\neg P_i$	0.2	1-x	1-y	z

- Given attributes and dependencies compute CPT
- Given networks but assuming some are hidden → GD methods

randomly initialized weights (CPT entries)

$P_w(b) = \prod_{d=1}^{|\mathcal{D}|} P_w(x_d) = \prod_{d=1}^{|\mathcal{D}|} \prod_{i=1}^{|X_i|} P(x_i | \text{parents}(x_i))$

in the summation $P(x_i = v_{ij} | \text{parents}(x_i) = u_{ik}) = w_{ijk}$

$w_{ijk} \leftarrow w_{ijk} - \alpha \frac{\partial \ln P_w(b)}{\partial w_{ijk}}$

$\sum_j w_{ijk} = 1$

- attributes but not dependencies → search through the model space
- no attribute nor dependencies → not doable

DIFFERENCES

- rules are by design exclusive and exhaustive
- rules do stand alone, even I need to follow from root to leaf

RULE BASED CLASSIFICATION → IF antecedent THEN consequent

$\text{Newness} = \# \text{ of tuples covered by a rule}$
 $\text{Newness} = \# \text{ of tuples correctly classified by a rule}$
 $\text{coverage}(R) = \frac{\text{Newness}}{|D|}$ $\text{accuracy}(R) = \frac{\text{Newness}}{\text{Newness}}$

CONFLICT RESOLUTION

- Size ordering → assign priority to the toughest rule with the most attribute
- Rule ordering → makes a list of rules and test it in order (implies the rejection of the rules that come before)
 - class-based ordering: decreasing order of prevalence or misclassification cost (no need to order inside a class)
 - rule-based ordering: by some quality measures (accuracy, cover, size) or after advice of domain expert

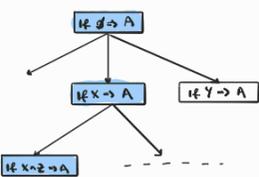
RULE MINING

- from DECISION TREE → one rule per leaf, they are mutually exclusive (no conflict) and exhaustive (one rule for all-val)
- pruned
 - condition that doesn't improve the accuracy of the rule can be pruned
 - any rule that doesn't improve the overall accuracy can be pruned

After pruning the rules will no longer be mutually exclusive and exhaustive

- I can class group the rules and find an order that minimize false positives → putting the next with lowest FP first
- The default class can be the one with most uncovered tuples

- HEURISTIC APPROACH → sequential covering algorithm: rules learned one at the time, remove the tuples covered, until no more training example or returned rule has quality under a specific threshold



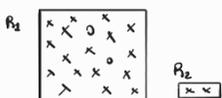
I use greedy depth-first search → from the empty one then consider each attributes I could test

$\text{FOIL-GAIN} = \text{pos}' \left(\log_2 \frac{\text{pos}'}{\text{pos}' + \text{neg}'} - \log_2 \frac{\text{pos}}{\text{pos} + \text{neg}} \right)$

To prune the rule I could use $\text{FOIL-PRUNE} = \frac{\text{pos} - \text{neg}}{\text{pos} + \text{neg}}$

- favour rules that have high accuracy and cover more positive tuples
- until $\text{foil-gain} > t \rightarrow$ add predicate
- increase with the accuracy of R on D-prune
- if the pruned version have higher foil \Rightarrow prune it
- conflict creation that needs avoiding

cannot use accuracy for rule quality, because I could have perfect accuracy, but very small coverage



$\text{pos}' =$ positive tuples that are covered by the new rule (correct)
 $\text{neg}' =$ negative tuples that are covered by the new rule (wrong) → new rule means with a predicate more

LAZY LEARNERS

→ less time training, more inference • use a richer hypothesis space that lead to better accuracy

- K-NN** → instance based learning
 - Voronoi diagram to visualize decision boundaries of 1-NN
 - with categorical values I can use the distance 0 if they are equal, 1 if they are different
 - if a value is missing from a tuple with all the others give the maximum distance
 - sorting and saving the training dataset to search-trees can reduce the # of comparison to $O(\log(n))$
 - can be parallelized
 - curse of dimensionality → less meaningful attributes must be eliminated
 - prediction with K-NN returning the mean of the neighbors, I can also weight by the distance

use locally weighted regression

- for discrete values K-NN return the most frequent value across the K nearest training example from \bar{x}_q
 $T(n, d, K) = O(n \cdot d \cdot K)$

1) WILSON EDITING → cleans overlap interclass regions → smoother boundaries
 $\forall o_i \in O$ label it with the K-NN and remove o_i misclassified

2) MULTI-EDIT → repeat Wilson to N random subsets
 partition O in $N \geq 3$ random subset, classify S_i with $S_{(i+1) \bmod N}$ discard incorrect edges, repeat until no changes

3) CITATION EDITING → $\forall o_i \in O$ find K-NN and C nearest cites (where o_i is in their K-NN)
 classify o_i both with K-NN and cites and discard the misclassified

4) SUPERVISED CLUSTERING → substitute O with O_c that are cluster representatives (one for cluster)
 - strong complexity and space reduction, but also accuracy

$$w_i = \frac{1}{d(x_i, x_j)^2}$$

compression rate = $(1 - \frac{c}{n}) \cdot 100$
 c = # of tuples after editing

- CASE-BASED REASONING** → solve new problems with past solutions to answer one
 - symbolic representations (tuples, graphs not points) where I use knowledge based reasoning (inference) eg patient X has Y symptom, treated with Z
 - when a new solution is found I will update my case-base to increase its competence → learning

ASSUMPTION: similar problems have similar solutions

- retrieve → need to define a similarity function. When |D| grows the efficiency of retrieval decrease: use particular data structures
- reuse → may need some adaptation
 - transformation using a function
 - generative approach (generate new from scratch)
- revise → get feedback as accuracy indicators
- retain → selective retention of new solutions, based on accuracy and the usefulness of saving it (since it will grow the dataset)

ENSEMBLE METHODS

→ use a combination of models in order to get a better accuracy

- BAGGING**: averaging the prediction over a set of n classifiers
 - I get n dataset from bootstrap sampling each of which allow me to learn a different classifier M_i
 - better accuracy since more detailed decision boundaries
 - computationally heavier, but highly preferable
 - by averaging I reduce the model variance make it more robust
- RANDOM FOREST** where the M_i are decision trees (> 100 trees)
 - create bi dataset using bootstrap, then at each node randomly choose m attributes from the available
 - the tree is fully grown and not pruned → low probability of overfitting (majority of votes)

$$E_T \left[\left(\hat{f}(x) - E_T[\hat{f}(x)] \right)^2 \right] \rightarrow \text{over different training sets from the same distribution } T \sim D$$

estimated median prediction on T_i - automatic estimation of attributes

- BOOSTING**: assign weights to each tuple. Then learn sequentially the classifiers updating every time the weights, in order to weight more the one misclassified by M_i (M_{i+1} can focus). M^* combines the outputs weighted by M_i accuracy
- better accuracy than bagging but also high risk of overfitting

ADABOOST initially $w_i = \frac{1}{n}$ At round K I get a bootstrap sample, learn $M_i \Rightarrow \text{err}(M_i) = \sum_{j=1}^n w_j \text{err}(x_j)$

- if a tuple is correctly classified I update the weight to $w_j' = w_j \cdot \frac{\text{err}(M_i)}{1 - \text{err}(M_i)}$, then $w_j'' = w_j' \cdot \frac{\sum w_j' \cdot \text{err}(M_i)}{\sum w_j'}$ → flip the same sum as before (-1)
- I assign a classifier weight as $w_i = \log\left(\frac{1 - \text{err}(M_i)}{\text{err}(M_i)}\right)$ → class with highest sum of weight vote wins
- reduce the model bias by sequentially focusing on solving errors

$$E_T[\hat{f}(x)] - f(x) \quad ; \text{variance-bias trade off}$$

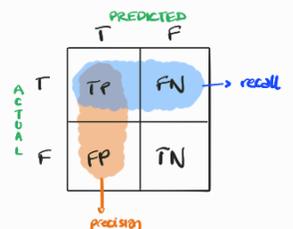
estimated by tree

MODEL EVALUATION

→ accuracy for unbalanced dataset doesn't work • I need to use the confusion matrix → I would like it diagonal

precision = $\frac{TP}{TP + FP}$ → % of tuples that the classifier labelled as positive are actually positive
 recall = $\frac{TP}{TP + FN}$ → % of positive tuples that the label classified as positive

- I can also create a multi-class confusion matrix by considering all the other attributes as negative



F-MEASURE $F_\beta = \frac{(1 + \beta^2) \cdot \text{precision} \cdot \text{recall}}{\beta^2 \cdot \text{precision} + \text{recall}}$

$\beta = 2$ weights recall twice as much

DATASET GENERATION for statistical evidence for comparing results of two models

1) **Hold-out method**: partition D in two randomly partitions K times, train K models, then calculate the metrics

2) **Cross-Validation**: partition dataset in K parts use D_i as test, all the others as train
 - maybe it's better to perform pre-processing only on training so I should avoid to do it to all dataset

- **LEAVE-ONE-OUT** K folds with $K = \#$ of tuples (for small size dataset) \rightarrow test only on one tuple
- **STRATIFIED CROSS-VALIDATION** \rightarrow maintain the same rate of majority/minority classes (better than the standard)

3) **BootStrap**: sample uniformly with replacement maintaining the original dataset cardinality $P(x \notin D_i) = (1 - \frac{1}{d})^d \approx \frac{1}{e} = 0.368$. get n dataset
 $Acc(M) = \frac{1}{K} \sum_i (0.632 \cdot Acc(M_i)_{test} + 0.368 Acc(M_i)_{training})$

\hookrightarrow if measured only on test is **negatively biased** since only .632 is seen \Rightarrow define this new estimator to correct the bias

M1 vs M2

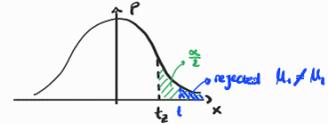
1) **T-TEST** (parametric test, assume two gaussian with similar variance) \rightarrow need to verify the assumptions before testing

- H_0 : the two distribution for $err(M_1)$ and $err(M_2)$ are the same (equivalent to $\mu_1 = \mu_2$ since $\sigma_1 = \sigma_2$ for hypothesis)
- $K-1$ degrees of freedom, with $K = \#$ of measurements

Use K folds and evaluate M_1 and M_2 on the same partitions

$$var(M_1 - M_2) = \frac{1}{K} \sum_{i=1}^K [err(M_{1i}) - err(M_{2i}) - (err(M_1) - err(M_2))]^2$$

$$t = \frac{err(M_1) - err(M_2)}{\sqrt{var(M_1 - M_2) / K}}$$



If the two models are compared on different K -folds: $var(M_1 - M_2) = \sqrt{\frac{var(M_1)}{K_1} + \frac{var(M_2)}{K_2}}$ $df = \min\{K_1, K_2\}$

• If $|t| > t_2$ we lies in the rejection zone $Z = \frac{|t|}{2} \rightarrow$ significance level because it is symmetric

2) **WILCOXON SIGNED RANK SUM TEST** (non-parametric, but variances must be similar)

H_0 : the medians of the two distribution are equals

$d_i = x_i - y_i \rightarrow$ then rank the $|d_i|$ with \pm for the minimum and so on

Label its rank with the sign of $d_i \rightarrow$ calculate W^+ and W^- ($W^+ + W^- = \frac{n(n+1)}{2}$)

• for H_0 the d_i will be centered around zero randomly: $W = \min(W^+, W^-) \rightarrow$ use a table of critical values to find P of getting a more extreme value

if $\frac{n(n+1)}{2} > 20 \rightarrow$ normal approximation $\mu_W = \frac{n(n+1)}{4}$ $\sigma_W = \sqrt{\frac{n(n+1)(n+1)}{24}}$

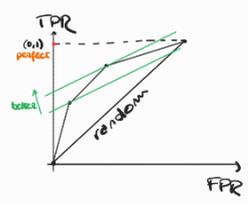
then I need to calculate the **z-score** $z = \frac{W - \mu_W}{\sigma_W}$ (how many std I'm outside the mean)

- if I have d_i that are equals \rightarrow give the same rank as the mean (eg 5.5) and reduce variance for each group of t tied ranks $\frac{t^3 - t}{12}$

- If $d_i = 0$ ignore that sample and adjust n accordingly

3) **AREA under the ROC CURVE (AUC)** \rightarrow alternative to F-measures for imbalanced dataset

$FPR = \frac{FP}{FP+TN}$ $TPR = \frac{TP}{TP+FN}$ \rightarrow area below ROC curve is a measure of the model accuracy independently from model thresholds \rightarrow single number



• the diagonal is like random choice because $TPR = FPR$

	predicted				
actual	<table border="1"> <tr> <td>0.1</td> <td>0.1</td> </tr> <tr> <td>0.3</td> <td>0.2</td> </tr> </table>	0.1	0.1	0.3	0.2
0.1	0.1				
0.3	0.2				

$cost = FPR \cdot P(neg) \cdot c(Y, n) + FNR \cdot P(pos) \cdot c(N, p) \Rightarrow 0 = dFPR \cdot P(neg) \cdot c(Y, n) - dTPR \cdot P(pos) \cdot c(N, p)$

$\rightarrow m = \frac{dTPR}{dFPR} = \frac{P(neg) \cdot c(Y, n)}{P(pos) \cdot c(N, p)} \rightarrow$ useful to get the min cost point of ROC and select right thresholds

• I can compute the AUC for each trial and then perform a statistic test on the result of AUC on the two models comparing single point ROC waves and comparing them statistically instead of using AUC of the full curve

IMBALANCED DATASET \rightarrow models will be biased towards the majority class

\rightarrow I need to rebalance just the training, not the test, otherwise the metrics will be biased

- 1) oversampling
- 2) under-sampling \rightarrow loose information
- 3) threshold moving - from 0.5, so that rare class tuples are easier to classify (less FN errors)
- 4) ensemble methods (boosting)
- 5) class weight, that modify the loss function

• **SMOTE (synthetic minority over-sampling technique)** \rightarrow create synthetic data instead of sampling (for minority class, typically improve AUC)

if point from the minority class compute the K -NN (from K points select the right number based on the oversampling required, like 2 out of 5 for a 20%)
 - then compute the difference between that point and one selected from K -NN and multiply by $\alpha \sim U(0,1)$ $P' = P + \alpha \cdot \text{difference}$

• pay attention to test it on original data



$\text{diff} = P'_{min} - P$

CLUSTERING

CLUSTER → set of objects that are similar within the cluster and dissimilar with other objects

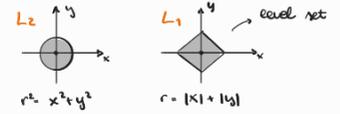
CLUSTER ANALYSIS → process of finding similarities between data and grouping similar data objects into clusters → **unsupervised learning**

Can be used as **preprocess** (data compression, outliers) or **stand-alone** tool

Similarity is expressed with a distance function $d(i,j)$ → change for the type of attribute
 • high intra-class sim. and low inter-class sim. can be **weighted** based on application data semantics → **weighted L2** I can get **ellipsoid** clusters

Quality of clustering → separate function that measure the quality. It's highly subjective

- single level vs. hierarchical partitioning → useful to have different clustering's level
- distance-based vs. connectivity-based similarity → euclidean or density
- exclusive vs. non-exclusive → sample can be associated with more clusters



Assessing clusters tendencies → evaluate if data set contains meaningful clusters

HOPKINS STATISTIC → spatial test to understand if a variable is randomly distributed across a space (uniformly distributed)

- 1) sample n points $\{p_1, \dots, p_n\} \subseteq D$ for each find its nearest neighbour in D $x_i = \min_{v \in D, v \neq p_i} \text{dist}(p_i, v)$
- 2) draw n points $\{q_1, \dots, q_n\} \sim U(D_{\text{max}}, D_{\text{min}})$ → same variation as D $y_i = \min_{v \in D} \text{dist}(q_i, v)$

→ $H = \frac{\sum y_i}{\sum y_i + \sum x_i}$ If uniform $H \approx 0.5$, otherwise $\sum x_i \ll \sum y_i \Rightarrow H \approx 1$

No: data is uniformly distributed

• compute H for several $\{p_1, \dots, p_n\}$ samples and average it
 If $H > 0.75 \Rightarrow$ cluster tendency with 90% of confidence

1) PARTITIONING METHODS → construct partition and evaluate them, then iterative approach to reassigning object from one partition to another

- I want to find K cluster such that they minimize the sum of squared distances with the cluster centroid/medoid $E = \sum_{i=1}^K \sum_{p \in C_i} (p - c_i)^2$ (centroid of C_i)
- I'm forcing the **compactness** of the clusters, hope to have **separation** as side effect (when $K = |D| \Rightarrow E = 0$)
- correct one this function to select K because it will always **decrease**

K-MEANS

c_i → centroid of the cluster $C_i = \frac{1}{|C_i|} \sum_{p \in C_i} p$

- 1) randomly assign K objects from D as centroid
- 2) reassign each object to the K clusters
- 3) recompute the centroid of the cluster

for each point n find the nearest cluster. $O(n)$ recompute cluster centers $O(n)$, repeat at most t iterations

$O(tkn)$ and since $t, k \ll n \Rightarrow$ linear with the number of objects

- unfortunately terminates at local optimum → multi-start approach
- need to specify K • not suitable for discovering non-convex shapes (only spherical)
- can be used only for continuous \mathbb{R}^n data (I would not know how find the center C_i) → I can use **K-modes** for categorical data

since clusters centers computed with the mean → L_2

categorically with highest frequency

- sensitive to outliers

→ I can use **K-modes** for categorical data

- change distance function for categorical
- we made to find row centroids (most frequent)

ELBOW METHODS → determine the optimal K and if there are clusters

Increasing K will decrease the intra-class variance of each cluster. I found the turning point when split cohesives clusters will not result in a big decrease $J = \sum_{i=1}^K \text{var}(C_i)$

K-MEDOIDS

medoid is the most centrally located object in a cluster

• can be used with more **d-type** and is more robust to outliers respect to K-means

- 1) randomly choose K representatives $\{o_1, \dots, o_k\}$
- 2) assign each object to the nearest representative o_i eg. with L_1 similarity
- 3) for each o_j randomly choose o_{random} and compute the cost of swap them when $S < 0 \Rightarrow$ replace them

$S = E'_{\text{new}} - E_{\text{old}} = \sum_{i=1}^k \sum_{p \in C_i} |p - o_i'| - \sum_{i=1}^k \sum_{p \in C_i} |p - o_i|$

can in another cluster

possible reassignment

PAM (original algorithm)

o_1 has the smallest $\sum_{i=1}^n d(i, o_1)$

$o_j, j > 1$ decreases the $\sum_{i=1}^n \min_{t=1, \dots, j} d(i, o_t)$ as much as possible

still local optimum because it cannot search through all possibilities (depends on initialization, I could be two swap from optimum)

BUILD PHASE

greedy choice (we do not explore every possibility)

SWAP PHASE

then consider all pairs $(i, j) : i \in \{o_1, \dots, o_k\} \wedge j \notin \{o_1, \dots, o_k\} \Rightarrow$ swap $i \rightarrow j$ that decrease the objective the most

non-medoid

• more robust to outliers than K-means

• doesn't scale well for large dataset $O(k(n-k)^2) \approx O(n^2)$

consider for each medoid the $n-k$ possible swaps and recompute for each S that cost $O(n)$
 $O(k(n-k)n) \approx O(k(n-k)^2)$ for each step

↳ the size will affect accuracy and complexity

CLARA (improvement from PAM) → apply on initial **sample** of the dataset and apply PAM

- if it is representative the medoids of the sample should be the one of the dataset
→ cannot find the best medoids if they are not on the sample
- to improve the approximation I can draw more samples and returning the best one

- 1) for m times draw a data sample and use PAM to get the k medoids
- 2) assign each object on the entire D to $\{0, \dots, m\}$
- 3) If the average dissimilarity on D is less than the current minimum → new model found

sample size → assign each $n-k$ remaining objects to the k medoids

$$O(m(KS^2 + K(n-K))) \rightarrow \text{can be used for larger dataset}$$

PAM $O(K(n-K)) \approx O(Kn^2)$

- I can introduce **bias** and accuracy loss for the sample

CLARANS (randomized clara) → clustering as searching a graph where each node represent the set of k medoids

- Two nodes are neighbour if they differ by just one medoid
- each node is associated with a **cost** that is the total dissimilarity → search a minimum
- At each step I will search all the neighbours and select the one with highest decrease → $K(n-k)$ connections per node, evaluate only m (constant)
↳ every medoid could be swapped by $n-k$ possibilities
- At each step I draw sample from the neighbours to examine → this is better than CLARA since I will not limit the search area
- I can search sequentially more local optimum, is a parameter of the algorithm
- more efficient and reliable than PAM/CLARA + better quality $O(mn) = O(n)$ → for each step compute m times S and swap

• All three alg can be seen as a graph search where PAM examines all the neighbours, CLARA draws an initial sample and search it all and CLARANS draw a random neighbours samples

2) HIERARCHICAL METHODS → create an hierarchical decomposition of the set

• since a step is taken it cannot be undone

- Use **distance matrix** as clustering criteria → no need for K but a **termination condition**
↳ **connectivity matrix**

- I need to define the distance between subset that can be constructed from the connectivity matrix, such distance is called **LINKAGE METRIC**
→ will affect the resulting clusters • defined as: $d(C_1, C_2) = \text{operation} \{ d(x, y) : x \in C_1, y \in C_2 \}$

- 1) minimum → **single link** (NN) distance of two clusters is the distance of their closest object → represent long chains } susceptible to outliers
- 2) maximum → **complete link** (furthest neighbour) : perform well with natural clumps, not with chains
• tends to minimize the diameter increasing of the cluster at each iteration
- 3) average → **pair-group average** : works well with clumps or elongated clusters

DENDROGRAM is a tree structure commonly used to represent the process of hierarchical clustering

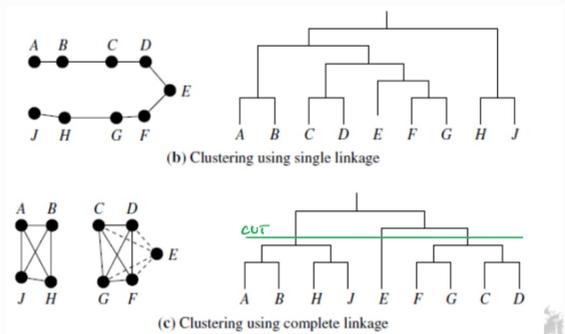
→ A cluster is obtained by cutting the dendrogram at the desired height
⇒ each connected component is a cluster

1) **single linkage algorithm** terminates when the minimum distance between nearest clusters is above t

If I see data as nodes ⇒ merge two clusters is adding an edge between the nearest objects
• the algorithm is also called **minimal spanning tree algorithm** because they are equivalent

2) **complete linkage** terminate when the maximum distance between nearest clusters exceed t →
each cluster is a **complete subgraph**, with edges all the node in the cluster

→ have at least $O(n^2)$ but in general gives better quality clusters



	A	B	C	D		
A	0	4	4	2		
B	4	0	3	4	AB	CD
C	4	3	0	4	AB	0
D	2	4	4	0	CD	2

$d=1$

AGNES (agglomerative nesting) → uses single link methods (merge)

- 1) Each point in a different cluster
- 2) with distance matrix found the nearest tuples and merge the clusters (simplify the dissimilarity matrix)

DIANA → uses average dissimilarity, inverse order of AGNES (split)
(diagnostic analysis)

• look the example on how to do it with dissimilarity matrix

1) start with a large cluster containing all the points

2) select the cluster with the highest intra-dissimilarity (distance between two points) → **largest diameter**

3) find the most disparate sample (largest average dissimilarity) → create a **splitter group** and populate the new cluster with points of old cluster that are closer to the splitter

$$V_i \text{ splitter } D_i = \text{avg } d(i, j) - \text{avg } d(i, j)$$

$\begin{matrix} j \in \text{splitter} & j \in \text{splitter} \end{matrix}$

- find the highest end if $D_i > 0$ ⇒ move it to the splitter group
- repeat until all $D_i \leq 0$ → recalculate since D_i will change (new point splitter)

4) repeat 2-3 until all clusters contain one point

BIRCH (balanced iterative reducing and clustering with hierarchies)

→ type of hierarchical clustering that proceed by merging

Agglomerative clustering for large amount of data → differently to other algorithms tries to minimize I/O because data can probably not fit in memory

$O(n)$

CLUSTERING FEATURES (CF) → summary of 0-th, 1-st and 2-nd moment of a cluster $CF = (N, LS, SS)$

N → # of points in the cluster

LS → linear sum of the points in the cluster $\sum_i x_i \in \mathbb{R}^p$

SS → squared sum $\sum_i |x_i|^2 \in \mathbb{R}$

I can use it to compute:

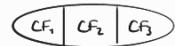
1) centroid of the cluster $\bar{x}_0 = \frac{\sum x_i}{N} = \frac{LS}{N}$

2) radius (avg. sq. root of squared distances from centroid) $R = \sqrt{\frac{\sum |x_i - \bar{x}_0|^2}{N}} = \sqrt{\frac{\sum |x_i|^2 - 2\bar{x}_0 \sum x_i + 2N\bar{x}_0^2}{N}} = \sqrt{\frac{SS - 2 \frac{LS \cdot LS}{N} + 2N \frac{LS \cdot LS}{N^2}}{N}} = \sqrt{\frac{N \cdot SS - |LS|^2}{N^2}}$

3) diameter $D = \sqrt{\frac{\sum_i \sum_j |x_i - x_j|^2}{n(n-1)}} = \sqrt{\frac{\sum_i \sum_j (|x_i|^2 + |x_j|^2 - 2x_i \cdot x_j)}{n(n-1)}} = \sqrt{\frac{2N \cdot SS - 2|LS|^2}{N(N-1)}}$

→ since they are all summations I can just update CF by adding the new elements or another CF when merging

CF-TREE is an **height balanced tree** \rightarrow hierarchical representation of clusters. Each node maintains the CF for each entry



I need the following parameters:

- 1) B (branching factor) → max # of children per non-leaf node → regulates the tree's height
- 2) T (threshold parameter) → maximum diameter of the clusters in the leaf node (no subcluster can have it higher)
- 3) L → max # of entries in a leaf

Phase 1: build an n -memory CF-tree incrementally. For each point that I consider:

- 1) If T violated (it will increase too much the diameter of the closest cluster)
 - a) If there is space → insert it as a new cluster in the leaf
 - b) if L violated ⇒ take the two further entries and create two nodes, put the remaining entries to the nearest leaf. update the entries and CF of the parent. we may split the parent as well (B violated). Splitting the root will increase the tree's height
- 2) otherwise insert it into the closest cluster (update CF)

→ If during tree construction CF tree exceed memory available I can reconstruct the partially-built tree in a smaller one from old tree leafs (no need to recompute points) (increase the value of T threshold). each leaf CF entry is treated as a **pseudo-point**

Phase 2 (optional): builds smaller CF-tree by increasing T from leafs

Phase 3: use an arbitrary cluster method for the CF-tree's leafs entries, that are treated as points (representative of **micro clusters**) → merge microclusters together. example I want to find K clusters

Phase 4 (optional): scan the entire DB to label point → I may need some outlier handling that doesn't fit in any cluster

- handles only numeric values and is sensitive to their order insertion
- cluster may not be so natural for the threshold given
- tends to be spherical (cannot exceed the given diameter)

CHAMELEON → measures similarity with a dynamic model

• two clusters are merged if the **interconnectivity** and **closeness** (proximity) are high relative to their **internal strength** (internal measures)

Construct the graph **K-NN** of the data → partition the graph into sub-clusters → **agglomerative clustering** of the sub-clusters (merge phase) → graph-partitioning algorithm

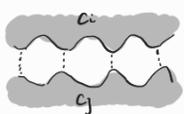
↳ agglomerative hierarchical clustering algorithm

- 1) K-NN captures the concept of neighborhood dynamically → important to set K correctly
- the radius of the neighborhood is inversely proportionate to the density of the region
 - the density is calculated from edge weight → each edge gets a weight based on how close the two nodes are
 - used to calculate inter-connectivity

2) Partition phase → a cluster is partitioned into subclusters C_i and C_j to minimize the sum of weights (edge cut) of the edges that would have been cut → absolute interconnectivity $EC_{\{C_i, C_j\}}$

- both $C_i \cap C_j$ must contain at least 25% of C nodes
- repeatedly select the largest sub-cluster (graph) and use the partition algorithm to bisect it until the largest contains fewer vertices than a specified number (15-5% of overall data)

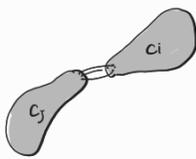
3) Agglomerative hierarchical clustering based on similarity that takes into consideration RI and RC



$$RI(C_i, C_j) = \frac{|EC_{\{C_i, C_j\}}|}{\frac{1}{2}(|EC_{C_i}| + |EC_{C_j}|)}$$

\downarrow relative interconnectivity
 - strength of the structure

\rightarrow absolute inter-connectivity (sum of weights of vertices between C_i and C_j)
 \downarrow internal inter-connectivity → minimum sum of weights of edges which would be cut by partitioning C_i into two roughly equal parts



$$RC(C_i, C_j) = \frac{\bar{S}_{EC_{\{C_i, C_j\}}}}{\frac{|C_i|}{|C_i| + |C_j|} \bar{S}_{EC_{C_i}} + \frac{|C_j|}{|C_i| + |C_j|} \bar{S}_{EC_{C_j}}}$$

\rightarrow avg weight between C_i and C_j
 \rightarrow avg weight of the edges that belongs to the min-cut bisector of C_j (same as for RI)

- measures the compactness

→ Merge the clusters that has RI and RC above TRI and TRC

- Visit all C_i and merge C_j with the highest absolute inter-connectivity $EC_{\{C_i, C_j\}}$ that satisfies $RI(C_i, C_j) \geq TRI \wedge RC(C_i, C_j) \geq TRC$

The time complexity is $O(n^2)$ in the worst case → arbitrarily shaped clusters with high accuracy

DENSITY-BASED METHODS → clustering based on densities

- discovers clusters of all shapes
- handles well noise
- one scan
- needs density parameters as termination conditions

Needs two parameters

- Eps: max radius of neighbour (hyper-sphere)
- MinPts: min # of points inside an ϵ -neighbour to consider it dense → the point is called CORE OBJECT

- directly density-reachable: an object p is ddr from q w.r.t ϵ, MinPts ↔
 - $p \in N_{\epsilon}(q)$
 - q is a core object $|N_{\epsilon}(q)| \geq \text{MinPts}$
- density-reachable: an object p is dr from q w.r.t ϵ, MinPts ↔ $\exists (p_1, \dots, p_n)$ with $p_1 = q, p_n = p$: p_{i+1} is ddr from p_i
 - not symmetric (if p not core object)
- density-connected: an object p is dc from q w.r.t ϵ, MinPts ↔ $\exists o$: p and q are dr from o
 - equivalence relation (a-a, transitive, symmetric)

DBSCAN → cluster is a maximum set of density-connected points

- no problem with the noise/outliers
- arbitrary shape $O(n^2)$

search the ϵ -neighbourhood of each point in the DB if above MinPts → create a new cluster and iteratively add the neighbourhood until no other point can be added

↓ otherwise noise

- noise → not core points (could be border region of a cluster)
 - outliers → noise that is not in any neighbourhood (not in cluster)
- the seq collects iteratively directly density-reachable points

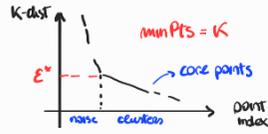
• Closure of density-connectedness to form clusters \rightarrow initially just a point, then applying the binary relation of density-connection I create the maximum closed set that is called cluster

• if I would use the dr for the cluster I would have missed points. If going in this sequential manner I would have had merge

A subset $C \subseteq D$ is a cluster \Leftrightarrow 1) internally connected $\forall o_1, o_2 \in C$ o_1 and o_2 are density connected
 2) maximality $\nexists o_1 \in (D \setminus C), o_2 \in C : o_1$ and o_2 are density connected

The parameters (ϵ and minPts) are valid for all the space, but I could have clusters with different density

• we define a function called **K-dist(p)** = distance with the k^{th} NN
 \rightarrow sort the point in descending order and plot it



OPTICS (ordering points to identify clustering structure)

• global density parameter cannot characterize the intrinsic cluster structure

\rightarrow I could use an hierarchical approach

- single link can merge two clusters just for a single connection chain
- dendrograms when I have a lot of points are difficult to understand

\rightarrow I could use density based algorithm but the # of parameters are infinite

\rightarrow if I discretize still the results would be hard to analyze

Density clusters are **monotonic** wrt their $\epsilon \rightarrow$ with fixed minPts higher density clusters (lower ϵ) are completely contained by lower density clusters (higher ϵ)

Given ϵ and minPts I can calculate the following properties of a point:

• **core distance**: smallest ϵ' such that p is a core-object \rightarrow distance from p to the minPts NN
 • if p is not a core-object \Rightarrow c-d is undefined

• **reachability distance** of object q from p is the minimum radius that makes q directly density-reachable from p
 $r-d = \max \{ c-d(p), dist(p, q) \}$ $\rightarrow q$ must be on its neighbourhood
 • can be undefined

not a real distance because not symmetric

\rightarrow or OrderSeed if not empty

1) Select a point from the input DB. Retrieve its ϵ -neighbourhood, calculate c-d and set its r-d to undefined

- if p is not a core object \rightarrow move with next element of OrderSeed (or DB if empty)
- if p is a core object $\rightarrow \forall q \in \epsilon\text{-neigh}(p)$ update their reachability distance from p and insert q into OrderSeed if not already processed

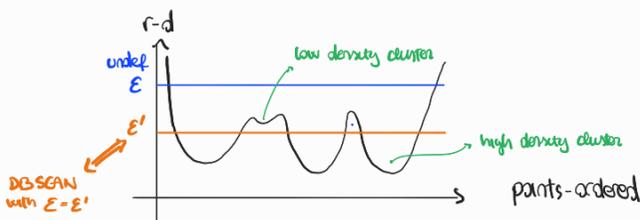
2) The object inside the **OrderSeed** are sorted by their reachability distance to the closest core obj

\rightarrow The next step of the while will retrieve the first element of OrderSeed

$O(n^2)$

\rightarrow once a point is examined is saved to a file with core-distance and its current minimum r-d

- The program output is a list of object (file) in the order the objects have been processed
- This is not a clustering algorithm since I don't assign objects but is just an ordering



\rightarrow I can get the clusters with different density parameters without re-executing the algorithm but just looking at the order

- wider valley correspond to more significant clusters
- peaks are threshold like noise

• with minPts = 2 \Leftrightarrow single link clustering
 - chaining effects problem with low minPts



DENCLUE → density algorithm can be highly sensitive to the value of radius used

- It is a clustering method based on a set of **density distribution functions**
- each observed object is treated as an indicator of high-probability density in the surroundings (depends on distance)

kernel density approximation of the probability density function: given $\{x_1, \dots, x_n\}$ iid from f random variable $\hat{f}_h(x) = \frac{1}{nh} \sum_{i=1}^n K\left(\frac{x-x_i}{h}\right)$

KERNEL: $K: \mathbb{R} \rightarrow \mathbb{R}^+$ must satisfy $\int_{-\infty}^{\infty} K(u) du = 1$ \wedge $K(-u) = K(u)$

an example $K\left(\frac{x-x_i}{h}\right) = \frac{1}{\sqrt{2\pi}} e^{-\frac{(x-x_i)^2}{2h^2}}$ or $K_{\text{square}}\left(\frac{x-x_i}{h}\right) = \begin{cases} 1 & \text{if } \left|\frac{x-x_i}{h}\right| \leq 1 \\ 0 & \text{otherwise} \end{cases}$

I can determine the clusters mathematically by identifying **clusters attractors** (local maxima of \hat{f})

- to avoid trivial local maxima I use a threshold where I consider x^* : $\hat{f}(x^*) \geq \epsilon$
- x^* are centers of clusters → to assign points to x^* I will use a **step-wise hill-climbing**

$$x^0 \leftarrow x \quad x^{j+1} \leftarrow x^j + \delta \frac{\nabla \hat{f}(x^j)}{|\nabla \hat{f}(x^j)|} \quad \nabla \hat{f}(x) = \frac{1}{nh} \sum K'\left(\frac{x-x_i}{h}\right) \cdot \frac{1}{h}$$

- it will stop when $\hat{f}(x^{k+1}) < \hat{f}(x^k)$ and assign x to attractor x^k
- I keep a list with points that have $d(x^i, x^j) \leq \frac{h}{2}$ and will assign these point to $x^* = x^k$ without re-applying the hill-climbing
↳ for each step of the climb $1 \leq j \leq k$

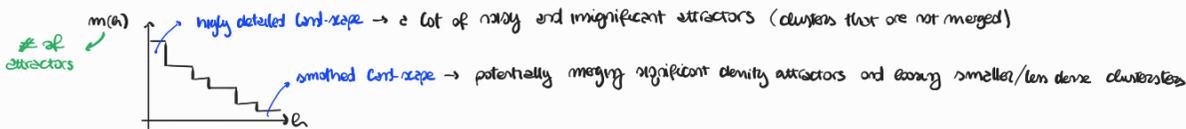
Noise or outliers are object x that converge to x^* : $\hat{f}(x^*) < \epsilon$

In order to get arbitrarily shaped clusters I need to merge density attractors

- A cluster $C \in \mathcal{D}$: $\forall x \in C, \exists x^* \in X$ ^{attractors set} $\hat{f}(x^*) \geq \epsilon$ and the rule below ($\forall x_i^*, x_j^* \in X$)
- I can merge two density attractors $\Leftrightarrow \exists$ a path $P \in F^d$ from x_1^* to x_2^* where $\forall p \in P: \hat{f}(p) \geq \epsilon$
↳ features space in d dimension

→ solid and compact math foundations. Extremely faster than existing density algorithms **$O(n \log n)$**

But need an accurate choice of h (determines the influence of the point in its neighbourhood) and ϵ (avoid noisy and insignificant attractors)



I should choose $h: [h_{\min}, h_{\max}]$ for constant $m(h)$ is larger → stable structure

• If the dB is noise-free all density attractors of \mathcal{D} are significant → $0 \leq \epsilon \leq \min_{x^* \in X} \{\hat{f}(x^*)\}$

grid computations to speed up the algorithm

1) The minimal bounding hyper-rectangle of dataset is divided into d -dimensional cubes, with edge length of 2ϵ

- numbered according to relative position from the origin: these keys will be stored in search-trees for easy retrieval

- for each populated cube (that has points inside) store key, # of points inside (N_c), pointers to those points and $\sum_{x \in C} x$
- connect neighbours cube $c_1, c_2 \in C_p \Leftrightarrow d(\text{mean}(c_1), \text{mean}(c_2)) < 4\epsilon \rightarrow$ if K is gaussian $\sigma = h$
↳ efficiently computed by the linear sum

2) Only highly populated cubes (C_{sp}) and the one connected to them are considered in determine clustering

$$C_{sp} = \{c \in C_p \mid N_c \geq \epsilon_c\} \quad C_r = C_{sp} \cup \{c \in C_p \mid \exists c_s \in C_{sp} \text{ connected to } c\}$$

$x \in C$ and $C_1, C_2 \in C$ we set $near(x) = \{x_i \in C_i \mid d(\text{mean}(C_i), x) < k\sigma \wedge \exists \text{ connection}(C_1, C_2)\}$

local density
 $\hat{f}(x) = \sum_{x_i \in near(x)} e^{-\frac{d(x_i, x)^2}{2\sigma^2}}$

\Rightarrow take the hill-climbing with \hat{f}

eg. \leftarrow allow to neglect marginal influences
 \rightarrow shouldn't I consider all distances to connected cubes?

CONSIDERATIONS FOR DENSITY-BASED CLUSTERING

- not effective when high dimension for clustering data \rightarrow sparse space
- unlikely that points are closer to each other rather than the avg distance of points

GRID-BASED CLUSTERING

\rightarrow use multi-resolution grid data structure, quantized space \rightarrow fast processing time but lose on Accuracy

STING (Statistical Information Grid)

- each cell is partitioned into a # of smaller cells \rightarrow at the lowest level are computed from dB, tested with a stat. test
- parameters of higher cells can be easily derived from smaller cell ($\mu, \sigma, \text{min}, \text{max}, \#, \text{type of distribution}$)
- top-down approach to answer spatial queries \rightarrow find regions with particular constraints or return attribute of a region
 \hookrightarrow I can start also from the non-top layer
- for each cell I compute the confidence interval reflecting the relevance of the cell in the query
 - calculated using the statistical parameters of each cell
 - label them as relevant or irrelevant \rightarrow then proceed to the lower level until the bottom is reached returning the relevant squares
- simple to parallelize, generation of clusters $O(n)$, query processing time $O(q)$ \rightarrow # of cell at the bottom
- Only vertical or horizontal boundaries of clusters \rightarrow no diagonal
- Is an approximation of region returned by DBSCAN as granularity $\rightarrow 0$

CLIQUE

\rightarrow identifies sub-spaces of the original high dimension space in which to perform a better clustering

- density-based + grid-based
- partitions each dimension by the same # of equal length intervals
- a unit is dense if the points in the unit are more than a threshold
- A cluster is a maximal set of connected dense units within a subspace

Identify dense units by the APRIORI principle

k-th dimensional unit is dense \Rightarrow its projection in (k-1)-dimension are dense.
 Therefore the candidates for being dense units are the ones generated by the dense (k-1) dense units

- 1) Identify in every dimension the dense intervals
- 2) Iteratively join two k-dimensional dense cells C_1, C_2 \rightarrow the join defined in this way avoid duplicates
 with $C_1 \subseteq \langle D_{i_1}, \dots, D_{i_k} \rangle$, $C_2 \subseteq \langle D_{i_1}, \dots, D_{i_{k-1}}, D_{j_k} \rangle$ \Rightarrow generate a (k+1) dimensional unit $\langle D_{i_1}, \dots, D_{i_k}, D_{j_k} \rangle$
share same k-1 dimensions and interval in those
- 3) Check if new has the necessary # of points specified by threshold \rightarrow and also if the other k-subspaces are dense?
 \rightarrow terminates when I cannot join anything or no dense results
- 4) Assembly dense cells in each subspace to find clusters \rightarrow find a convenient representation
 - Minimum description length principle (MDL) \rightarrow use maximal regions to cover connected dense cells
 \hookrightarrow hyper-rectangle containing only dense cells, and cannot be expanded in any dimension
 - the best description of a cluster is NP-hard, CLIQUE use a greedy approach, where it starts from a cell and then try to expand the rectangle in any dimension. Then add a new rectangle
 - \rightarrow terminates when all the dense cells are covered

- scales linearly with the # of data in input and also with the dimensions of the data
- needs proper tuning of grid size and density threshold → I get get low accuracy for the simplicity of the method
- automatically find sub-spaces of highest dimensions in which the clusters exist (the projections)

EVALUATING CLUSTERS → emerging cluster tendencies ($x \notin U$), # of clusters and clustering quality

CROSS VALIDATION METHOD → use a part to perform evaluation (eg. sum of squared distances from nearest centroid)
 $\forall k > 0$ that is interesting I can take the avg of the m results and comparing the overall quality

CLUSTERING QUALITY MEASURES

1) **EXTRINSIC** → supervised $Q(C, C_g)$ → quality of cluster C compared to the ground-truth C_g

must satisfy:

- homogeneity → the purer the better
- completeness → assign all equal class obj to the same C_i
- Rag Bag** → putting wrong obj to a pure cluster should be penalized more than putting inside a miscellaneous cluster
- Small cluster preservation: splitting a small category to pieces is more harmful than splitting a larger one

BCUBED → evaluates precision and recall for every object

precision → indicates the **purity** of that category inside my cluster
 recall → indicates the **completeness** of the category into the real cluster

Let $D = \{o_1, \dots, o_n\}$ a set of object and C a cluster on D

Let $L(o_i)$ the object category and $C(o_i)$ what the clustering assigned (a cluster id)

$$\Rightarrow \text{correctness}(o_i, o_j) = \begin{cases} 1 & \text{if } L(o_i) = L(o_j) \Leftrightarrow C(o_i) = C(o_j) \\ 0 & \text{otherwise} \end{cases}$$

Precision = $\frac{1}{n} \sum_{i=1}^n \frac{\sum_{o_j: i \neq j, C(o_i) = C(o_j)} \text{correctness}(o_i, o_j)}{|\{o_j: i \neq j, C(o_i) = C(o_j)\}|}$ → # of correct clustering inside of o_i 's cluster

-hope they tends to one

Recall = $\frac{1}{n} \sum_{i=1}^n \frac{\sum_{o_j: i \neq j, L(o_i) = L(o_j)} \text{correctness}(o_i, o_j)}{|\{o_j: i \neq j, L(o_i) = L(o_j)\}|}$ → # of correct clustering inside o_i real cluster

2) **INTRINSIC** when ground true is not available → evaluate clustering based on compactness and separation

SILHOUETTE coefficient

$\forall o \in D$ compute:

- $a(o)$ = avg distance between obj and it's cluster obj
- $b(o)$ = min avg distance to a cluster where o doesn't belong to

low better
 $a(o) = \frac{\sum_{o' \in C_i, o' \neq o} \text{dist}(o, o')}{|C_i| - 1}$

high better
 $b(o) = \min_{C_j: j \neq i} \left\{ \frac{\sum_{o' \in C_j} \text{dist}(o, o')}{|C_j|} \right\}$

$S(o) = \frac{b(o) - a(o)}{\max\{a(o), b(o)\}}$

preferable
 $S(o) \in [-1, +1]$

when $S(o) < 0$ → a different cluster is in avg nearer → should be avoided

- then I can take the avg of all points
- usable \Leftrightarrow the clusters are **convex**

MINING FREQUENT PATTERN

- PATTERN:**
- a set of items {milk, bread}
 - a subsequence (PC, monitor) → the order matters
 - a substructure pattern → like a subgraph or a sub-tree

FREQUENT PATTERN is a pattern that occurs frequently in a data set
 → finding them is essential for mining associations and find correlations

Itemset $I = \{I_1, \dots, I_m\}$ **k-itemset** $X = \{I_1, \dots, I_k\}$

Absolute Support of X → frequency of occurrence of an itemset (how many times it shows up in D)

Relative Support $s =$ fraction of transactions that contain X → probability of X in being in a transaction

• An itemset X is frequent when its support is bigger than a threshold

Transaction T is a set of items: $T \subseteq I$

• an **association rule** is an implication $X \rightarrow Y$ where $X \subseteq I, Y \subseteq I, X \cap Y = \emptyset$

• **Support** of $X \rightarrow Y$ on D (transaction DB) → $\text{supp}(X \rightarrow Y) = P(X \cup Y)$ → probability that a $T \in D$ is $T \supseteq X \cup Y$
 - reflects usefulness of the rule (in 2% of transactions are purchased together)

• **conf** ($X \rightarrow Y$) = $P(Y|X) = \frac{P(X \cap Y)}{P(X)} = \frac{\text{supp}(X \cup Y)}{\text{supp}(X)}$
 - reflects certainty on the rule (70% of who bought a PC also bought a TV)

I can associate each item a boolean variable and create a **boolean vector** to keep track of what do I have in a transaction (not how many)

In general to mine rules:

- 1) Identify all frequent itemset (at least min-supp occurrences)
 - if an item is not frequent, useless to add more to it because will be more infrequent
- 2) Generate strong association rules that satisfy minsup and minconf

{Beer}: 3, {Diaper}: 4, {Beer, Diaper}: 2

{Beer} → {Diaper}, {Diaper} → {Beer}
 → same support, different confidence

$\{a_1, \dots, a_n\}$ contains $\sum_{k=1}^n \binom{n}{k} = \sum_{k=1}^n \frac{n!}{k!(n-k)!} = 2^n - 1$ sub-itemsets
 → mine closed itemset and max-itemset instead

X is **closed** in a dataset D if X is frequent and \nexists a super-itemset $Y \supset X$ with same support as X
 X is a **max-itemset** in D if X is frequent and \nexists a super-itemset $Y \supset X$ that is frequent

C → set of all closed FI in D , satisfying minsup

• contains **complete information** to create the whole set of frequent itemsets (C and its count information) → Data Compression

M → set of all max-frequent-itemset in D , satisfying minsup

↳ one all subsets have the same count

- in the worst case M ^{max len transaction} frequent item-set can be generated
 ↳ # distinct elements

- depends strongly on minsup threshold

A-PRIORI ALGORITHM

downward closure property → any non-empty subset of a frequent itemset must be frequent

- anti-monotonicity property: If C fails a test also all its supersets will fail

A-priori principle: If an itemset is infrequent \Rightarrow I should not test any of its supersets

1. scan to get the frequent 1-itemset
2. get (k+1)-itemsets from combining k-itemsets
3. test candidate on D

↳ one DB test after each C_k

frequent itemset of length k

How to construct L_k from L_{k-1} with the a-priori principle

1. join step $C_k \leftarrow L_{k-1} \bowtie L_{k-1}$

- assuming that transactions and itemsets are sorted in lexicographic order
- $l_i[k-2]$ refers to the $(k-2)^{th}$ term in the l_i itemset
- $l_i[1] < l_i[2] < \dots < l_i[k-1]$ $l_i \rightarrow$ is a $(k-1)$ -itemset

\rightarrow since they are ordered $L_{k-1} \bowtie L_{k-1}$ is joinable on two sets $l_1, l_2 \Leftrightarrow$ the first $k-2$ elements are equal (since all the subsets must be present in L_{k-1} , then also the one that has the first $k-2$ elements in common) and $l_1[k-1] < l_2[k-1]$ \rightarrow avoid duplication

2. prune step check that each $(k-1)$ -subset of a candidate k -itemset is present in L_{k-1} done for $k > 2$ (1-itemset are frequent)

In order to generate association rules I can:

$\forall e \in \text{frequent_itemset}$ $\exists s \subset e, s \neq \emptyset$ non empty subset
 output $s \Rightarrow e-s \Leftrightarrow \frac{\text{supp_count}(e)}{\text{supp_count}(s)} \geq \text{min_conf}$

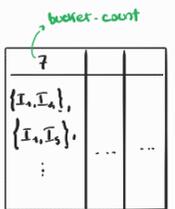
RARE ITEM PROBLEM \rightarrow items that occurs really infrequently and are pruned, but would generate interesting rules
 happens in strongly unbalanced dataset

- e problem with confidence is that is sensitive to the consequent (what comes after the arrow)
 \rightarrow consequent with higher support will automatically produce higher confidence even though no association exists

IMPROVEMENT of the APRIORI

1) If I partition $D \Rightarrow$ an itemset must be frequent in at least one of the DB equal partitions
 $\text{supp_count}(x) \geq \text{min_sup} \cdot |D|$

- I can find the frequent itemsets in the local partitions in parallel
- Combine all local frequent itemsets to form candidate itemsets \rightarrow check it scanning the DB
- just two scan since all the local partitions fits in memory



2) Hash-Boxed-Technique to reduce the size of the candidate C_k . When scanning for k -itemsets I can create an hash-table in which I map all $(k-1)$ -itemsets and associate for each bucket a counter, not consider the one with counter lower than a threshold

- Generation of candidate itemsets is expensive \rightarrow a lot of possible combinations (exponentially to the max-length of the frequent itemset)
- generate-and-test approach \rightarrow BREADTH-FIRST SEARCH
- A lot of DB scan (I/O operations)

ECLAT \rightarrow exploring vertical data format (equivalence class transformation)

$\{TID : \text{itemset}\} \rightarrow$ horizontal format
transaction id

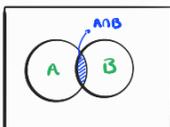
$\{I1, I7, I3\}$
 $t(AB) = \{T1, T2, \dots\}$
 $\{ \text{item} : TID_set \} \rightarrow$ vertical format
set of transaction ids

$t(x) = t(y)$ always happen together $t(x) \subset t(y)$ transaction having x always has y

- Use of a-priori property to generate $(k+1)$ -itemsets
- No need to scan the DB to find k -itemset support ($k > 1$). However TID list can be quite long and expensive to intersect

EVALUATION

$P(\text{videos}) = 75\%$ r : computer games \rightarrow videos $\text{sup}(r) = 40\%$ $\text{conf}(r) = 66\%$
 \rightarrow negatively correlated, buying PC games decreases the probability of buying videos



\rightarrow from association analysis to correlation analysis
 CORRELATION LIFT

$\text{lift}(A \Rightarrow B) = \text{lift}(B \Rightarrow A) = \frac{\text{conf}(A \rightarrow B)}{\text{sup}(B)} = \frac{\text{conf}(B \rightarrow A)}{\text{sup}(A)} = \frac{P(A \text{ and } B)}{P(A)P(B)}$

\rightarrow use solve the problem if B is always present

- correlation is symmetric instead of association rule is directional
- Lift measures how many times more often X and Y occur together respect than expected
 If lift < 1 → negatively correlated, the presence of one reduce the other
 → it is not down-ward closed ^{it's not monotonic, like sup-count, because probabilities can change in unpredictable ways} and doesn't suffer from the **rare item problem**
- possible problem when two rare items occur together by chance and produce an enormous lift value

Another metrics that I could use is the χ^2

$$\text{all-conf}(A,B) = \frac{\text{sup}(A \cup B)}{\max\{\text{sup}(A), \text{sup}(B)\}} = \min\{P(A|B), P(B|A)\}$$

$$\text{max-conf}(A,B) = \max\{P(A|B), P(B|A)\}$$

$$\text{Kulczynski}(A,B) = \frac{1}{2} [P(A|B) + P(B|A)]$$

→ can be 0.5 for complete irrelevance or just one side is important (dataset imbalance) → I need another number

0.5 for imbalance or —
 $E[0,1]$ closer to 1 stronger the relation
 null-invariant

A null transaction is a transaction that does not contain any of the item-sets being examined

- lift and χ are not null-invariant, since adding null transactions increase the lift value

$$\frac{\frac{n_{ano}}{n_r+x}}{\frac{n_a}{n_r} \cdot \frac{n_b}{n_r+x}} = \frac{n_{ano}}{n_a n_b} \cdot (n_r+x) = \text{lift-old} + x \frac{n_{ano}}{n_a n_b}$$

• important property especially in big dataset

We can measure the imbalance in two itemset A and B

$$\text{IR}(A,B) = \frac{|\text{sup}(A) - \text{sup}(B)|}{\text{sup}(A) + \text{sup}(B) - \text{sup}(A \cup B)}$$

Tends to one when completely imbalanced (eg $\text{sup}(B) \rightarrow 0 \rightarrow \text{sup}(A \cup B) \rightarrow 0$)

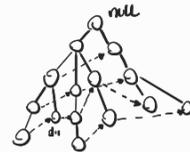
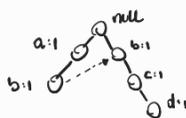
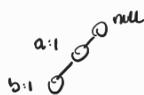
- generally use IR and Kulc to get a good picture of a rule

FP-GROWTH → avoid candidate generation that is the apriori bottleneck

- builds a compact data structure of \mathcal{D} that fits into memory (in two scan) and then find the FI

1) FP-TREE construction

- 1) scan for each single item, discard the infrequent, then sort elements in a transaction by decreasing order, based on their support → more overlapping (heuristic, not always true)
- 2) for each transaction, map it to a path in the tree: each node is an item and has an associated counter
 - if T1 and T2 has the same prefix, they will have overlapping path → **DATA COMPRESSION**
 - pointers are maintained between nodes containing the same items



independent of T order

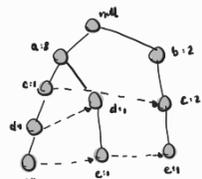
- **best case**: 1 path FP-tree

- **worst case**: every transaction doesn't have common prefixes
 • heavier than the original \mathcal{D} (need pointers and counters)

2) FI generation

- **bottom-up** generation (from leaf to root) → divide and conquer approach

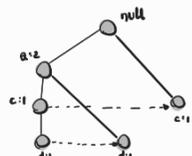
sub tree containing e:
 → FI ending in e



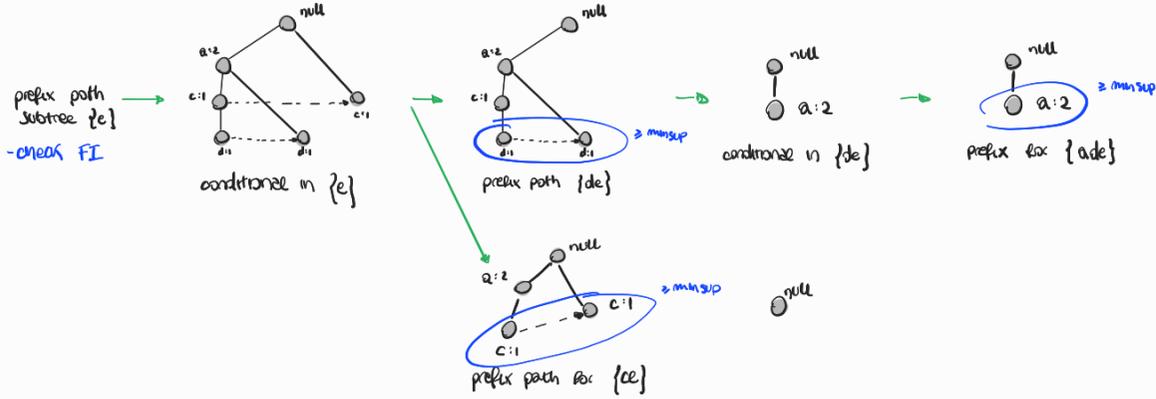
- process recursively the trees to extract FI that ends with e - de - cde side
 • **divide and conquer**

conditional

sub-tree containing e: FP-Tree that contains only T with e (then remove it)
 - remove items that cannot be frequent - update counters



1. Obtain the prefix path sub-tree for e
2. Check if e is frequent (adding counts along the linked list) → if yes add {e} to FI
3. If {e} ∈ FI → build the e conditional tree and check for {ce}, ~~{ae}~~, {de} and {ae} recursively



$$FI_1 = \{ \{de\}, \{ade\} \}$$

$$FI_2 = \{ \{ce\} \}$$

- FP tree is heavy to build, but I compress D (could discard it and save just FP if it doesn't grow)
 - **completeness** (preserve complete info)
 - **compactness** (reduce irrelevant info + common path)
- Order of magnitude faster than the apriori, especially for small min_sup

APRIORI → BFS, generate and test candidate itemset, level by level

FP-GROWTH → DFS, explore one branch of FP-tree completely before moving on

ADVANCED FREQUENT PATTERN ANALYSIS

- in **multi-level association rule** I should reduce minsup accordingly for lower level item (say milk)

$$\begin{matrix} \text{milk} \rightarrow \text{bread} \\ \text{say milk} \rightarrow \text{bread} \end{matrix} \quad \frac{\text{supp}(\text{milk})}{\text{supp}(\text{say milk})} \approx \frac{\text{conf}(r_1)}{\text{conf}(r_2)}, \quad \text{conf}(r_1) \approx \text{conf}(r_2) \quad \text{then } r_2 \text{ is redundant because of the ancestor}$$

- I should use non-uniform, group based min sup (more valuable items are less frequent)

- multi dimensional rules (≥ 2 predicates)

- interdimensional rule: $\text{age}(x, "19-25") \wedge \text{occupation}(x, "student") \rightarrow \text{buys}(x, "coke") \rightarrow$ no repetition
- hybrid dimensional rule: $\text{age}(x, "19-25") \wedge \text{buys}(x, "popcorn") \rightarrow \text{buys}(x, "coke")$

- for quantitative attributes I can discretize / cluster

NINE NEGATIVELY CORRELATED PATTERNS

- negatively correlated patterns infrequent tends to be more interesting than frequent one

def1: If X and Y are both frequent but rarely occurs together $P(X \cup Y) < \text{sup}(X) \cdot \text{sup}(Y)$
 \rightarrow X and Y are negatively correlated. *it's not null invariant (depends on the # of null transaction)*

def2: $\text{sup}(X \cup Y) \cdot \text{sup}(\bar{X} \cup \bar{Y}) \gg \text{sup}(X \cup Y) \cdot \text{sup}(\bar{X} \cup \bar{Y})$ - still suffer from null invariant

def3: X and Y frequent $\frac{1}{2}[P(X|Y) + P(Y|X)] < \epsilon \rightarrow$ are negatively correlated

CONSTRAINED BASES MINING

- knowledge constraint: classification we want the posterior to be the class label
- data constraint: products sold in Chicago
- dimension constraint: filter for relevance to a specific hierarchy of data (eg. say milk)
- pattern constraint: specify the form of conditions to be mined (predicates, # of conditions, ...) \rightarrow metarules

META-RULES: form an hypothesis regarding the relations the user is trying to confirm

$$P_1 \wedge P_2 \dots \wedge P_n \Rightarrow Q_1 \wedge \dots \wedge Q_m$$

I can use the rule constraints to prune the search space:

- pruning **pattern** search space (not consider some candidate patterns in the generate and test phase)
 - anti-monotonic: if constraint c is violated, it's further mining is recomputed (it's superset will fail) $\text{sum}(I.\text{price}) \leq 100 \$$
 - monotonic: if c is satisfied, no need to check c again (supersets will hold) $\text{sum}(I.\text{price}) \geq 100 \$$
 - succinct: c must be satisfied, so I can start with dataset where c holds \rightarrow avoid to enumerate everything $\text{min}(J.\text{price}) \geq 50 \$$ \rightarrow avoid generate and test (no support counting)
 - convertible: c not monotonic or anti-monotonic, but it can be converted into if items in the transaction are ordered \rightarrow any price $\geq x$
 - strongly convertible when monotonic or anti-monotonic according to the order
 - may exist **conflicts** on the order of items required \rightarrow prune with one, then the other
- pruning **data** search space (consider only useful part of the dataset) \rightarrow remove transaction
 - data succinct: remove some transaction before pattern mining
 - data anti-monotonic: if a transaction t doesn't satisfy c, t can be pruned from further mining

MINING COLossal FREQUENT PATTERN \rightarrow itemset of cardinality ≥ 50

- curse of downward closure property \rightarrow any subpattern of a frequent pattern is frequent \rightarrow the number of subitemset that are frequent are exponential
- considering closed itemsets may partially alleviate the problem but not really solve it, also them could be exponential
- closed patterns are usually attached with greater importance than those of small pattern and only a small # of patterns are closed
- \rightarrow Apriori or FP-Growth will get stuck in the huge number of **middle sized patterns** \hookrightarrow stuck in huge # of subtrees

PATTERN-FUSION traverse the tree in a **bounded-breadth way** \rightarrow approximation that quickly find most of them

1. only a subset of patterns are used as starting points \rightarrow avoid exponential search space
2. at each step it agglomerates more items together \rightarrow shortcut to reach massive size

Core Patterns \rightarrow for a frequent pattern α , a **subpattern** β is T-core pattern of $\alpha \Leftrightarrow \beta$ shares a similar support set

$$\frac{|D_{\alpha}|}{|D_{\beta}|} \geq r \quad r \in (0, 1] \quad |D_{\alpha}| \rightarrow \# \text{ of patterns containing } \alpha$$

core ratio

A pattern α is (d, T) -robust if d is the max # of items that can be removed from α for the remaining pattern to stay a T -core pattern of α

$$d = \max_{\beta} \{ |\alpha| - |\beta| : \beta \subseteq \alpha \wedge \beta \text{ is a } T\text{-core of } \alpha \}$$

$d \in \mathbb{Z}(2^d)$

maximize patterns have for more core-patterns than smaller ones \rightarrow remove some items and maintain similar support $\text{Dist}(\alpha, \beta) = 1 - \frac{|D_\alpha \cap D_\beta|}{|D_\alpha \cup D_\beta|}$

If two pattern α, β are core patterns of the same pattern \rightarrow they have their distance bounded by a ball $\text{Dist}(\alpha, \beta) \leq 1 - \frac{1}{\frac{2}{T} - 1}$

- given one pattern from its core pattern subset I can get all of them, just by examining the bounding ball
- If maximize pattern, the ball will be denser

If we draw randomly from the set of patterns of size c we would be more likely to pick a core-dependent of a colossal pattern than of a small pattern
 - merging more core patterns to get the maximize one

1) Initial Pool \rightarrow use an existing algorithm to mine all frequent pattern up to a small size

2) Iterative Pattern Fusion

- K seeds pattern are randomly picked from the current pattern pool \rightarrow high probability of being core-dependent of a colossal pattern
- for each one I get those within a bounding ball centered in the pattern (from the current pool)
 - \rightarrow all the pattern found are merged together to generate a set of super-patterns. This set will be the new pool for the next iteration
 - \rightarrow avoids mid-sized patterns, more than one for each K
- If the pool exceed a threshold I will sample it
- check that the set returned is frequent?

3) Termination when the current pool contains no more than K patterns at the beginning

\rightarrow The bigger the pattern, the greater the chance to be generated

MINING COMPRESSED PATTERNS

$$D(P, Q) = 1 - \frac{|I(P) \cap I(Q)|}{|I(P) \cup I(Q)|} \rightarrow \text{distance between two patterns}$$

Given two pattern A and B, A can be expressed by B if $O(B) \subseteq O(A)$
 - where $O(A)$ is the itemset of pattern A

$$\text{In a cluster } c: \bigcup_{i=1}^k O(P_i) \subseteq O(c)$$

A pattern P is δ -covered ($0 \leq \delta \leq 1$) by pattern P' if $O(P) \subseteq O(P') \wedge D(P, P') \leq \delta$
 - could belong to multiple clusters

Finding the minimum set of representative pattern is NP-hard

REDUNDANCY AWARE TOP-K PATTERNS \rightarrow not only high frequency, but also small redundancies

$S(p)$ is the degree of usefulness of the pattern P \rightarrow can be based on subjective or objective measures (supp, conf, ...)

\downarrow significance

$$R(p, q) = S(p) + S(q) - S(p, q) \rightarrow \text{we can approximate it with distance between patterns}$$

\downarrow redundancy

I want to find the K patterns that maximise the marginal significance \rightarrow relevant and not redundant to the already picked

OUTLIER ANALYSIS

outlier → object that deviates significantly from the normal object as if they were generated by a different mechanism
 - unusual credit card purchase - anomaly detection is different (weas points are merged in the model)

noise → random error in measuring
 - should be removed before outlier detection → may hide outliers

GLOBAL OUTLIER → points significantly deviates from the rest of the data

CONTEXTUAL OUTLIER → deviates significantly based on a selected context (from contextual attributes) → based on behavioural attributes (characteristics of the object)
 - 35°C in summer or winter

COLLECTIVE OUTLIERS → subset of data collectively deviate
 - intrusion detection in IT systems (multiple bots)

• I could have multiple types of outliers, one object may belong to more than one type
 • domain specific distance measures

1. Based on user labelled outlier (supervised, semi-supervised, unsupervised)

• supervised → classification problem

- imbalanced classes → boost outliers
 - recall is more important (not modelling outliers for normal data) → catch as many outliers as possible

• unsupervised → anomaly objects or divided into multiple groups. A cluster is distinct from all groups

• cannot detect collective outliers effectively • high FP rate and still miss some outlier
 • could we clustering algorithm directly → no cluster = outliers, but difficult from noise and costly since I need to cluster everything

• semi-supervised → when the # of labelled data is small

- depends if I have label for normal or outliers - use proximity measures and unsupervised help

2. Based on assumptions about normal data (statistical, proximity, clustering)

• STATISTICAL TECHNIQUES → effectiveness largely depends on whether the assumption of the statistical model holds or not
 model-based methods (assume underlying distributions)

Assumes data are generated by a stochastic process → generative model

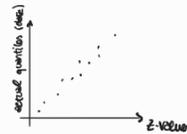
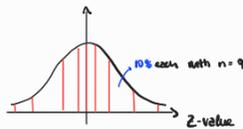
• learn a generative model and identify objects in low probability regions as outliers

parametric → data generated by a parametric distribution $f(x, \theta)$ ^{fixed in advance} the smaller this value, the more likely x is an outlier
 ↳ approximate

NORMAL DISTRIBUTION → assume univariate data (one attribute)

1. Using Q-Q plots → useful in general to understand if a sample comes from a continuous distribution

- plot quantiles from the sample against the quantile of the normal distribution
 • divide the N into $n+1$ segments, where n is the number of values → then calculate the z-values (cut-off for each segment) that is the # of std from the mean
 • if straight line → data comes from normal (the slope depends on the σ of the data)



2. Maximum Likelihood to estimate μ, σ

$$P(x_i | \mu, \sigma^2) = \frac{1}{\sqrt{2\pi\sigma^2}} e^{-\frac{(x_i - \mu)^2}{2\sigma^2}} \rightarrow P(X | \mu, \sigma^2) = \prod_{i=1}^n P(x_i | \mu, \sigma^2)$$

Find μ, σ^2 maximizing the likelihood $N(\mu_0, \sigma_0^2) = \operatorname{argmax} \{ P(X | \mu, \sigma^2), \mu, \sigma^2 \in \mathbb{R} \}$

$$\ln P(X | \mu, \sigma^2) = \sum_{i=1}^n \ln P(x_i | \mu, \sigma^2) = -\frac{n}{2} \ln(2\pi) - \frac{n}{2} \ln \sigma^2 - \frac{1}{2\sigma^2} \sum (x_i - \mu)^2$$

$$\frac{\partial \ln P(X | \mu, \sigma^2)}{\partial \mu} = -\frac{1}{\sigma^2} \sum (x_i - \hat{\mu}) = 0 \Rightarrow \hat{\mu} = \frac{1}{n} \sum x_i$$

$$\frac{\partial \ln P(X | \mu, \sigma^2)}{\partial \sigma^2} = -\frac{n}{2\sigma^2} + \frac{1}{\sigma^4} \sum (x_i - \mu)^2 = 0 \Rightarrow \frac{1}{\sigma^2} \left[\sum (x_i - \mu)^2 - n\sigma^2 \right] = 0 \Rightarrow \sigma^2 = \frac{1}{n} \sum (x_i - \mu)^2$$

$\mu \pm 3\sigma$ contains 99.7% of data, everything outside can be considered an outliers

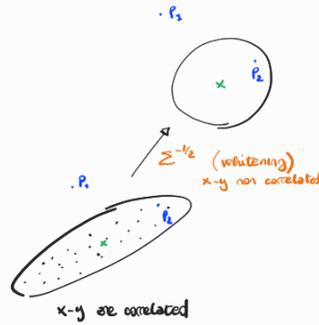
3. Grubb's Test → find whether a min/max value is an outlier

- find **G test statistics** order and find mean (\bar{y}) and std (s) of dataset

$$G = \frac{\max_{i=1, \dots, n} |y_i - \bar{y}|}{s}$$
 → this is because only max or min
- find G critical value from a table (given α and N)
- if $G_{test} < G_{critical}$ keep the sample, not an outlier, otherwise reject

MULTIVARIATE DATA → transform multivariate outlier decision to univariate
 PARAMETRIC METHODS

1. Mahalanobis distance: between a point and a distribution
 - euclidean distance works well when dimension are equally weighted and independent
 - The steps are:
 - Transform the columns into uncorrelated variables
 - scale to have variance $\sigma^2 = 1$
 - calculate the euclidean distance



$$MDist^2(0, \bar{0}) = (0 - \bar{0})^T \Sigma^{-1} (0 - \bar{0})$$

← covariance matrix → use in z score divide by σ

→ If the variables are strongly correlated the covariance will be high and the distance reduced

Use the distance to perform univariate outlier analysis → Grubb's test
 • computationally heavy since I need to find S and S^{-1} , also store them

2. χ^2 statistic
$$\chi^2 = \sum_{i=1}^n \frac{(x_i - \bar{x}_i)^2}{\sigma_i^2}$$
 $n = \# \text{ of variables}$

- when n is big enough (> 30) → $\chi^2 \sim N$ (central limit theorem)
- when χ is really large → outlier
- not the upper limit as: $\bar{x}^2 + 3Sx^2$

NON PARAMETRIC METHODS → not assume any a-priori statistical model and determine the model from the input data (not only the parameters)
 - parameters are flexible and learned by the data, not fixed in advance

1. Histogram → hard to choose appropriate bin size:
 - too large then outliers in frequent bin (false negative)
 - too small then normal objects in rare bins (false positive)

→ better to use kernel methods

2. Kernel Density Estimation
$$\hat{f}_n(x) = \frac{1}{n} \sum_{i=1}^n k_h(x - x_i) = \frac{1}{nh} \sum_{i=1}^n K\left(\frac{x - x_i}{h}\right)$$

 ← symmetric, not necessarily ≥ 0 , but integrates to one
 ← smoothing parameter

The complexity of statistical methods depends on the model

- gaussian → linear
- for kernel density estimation, the model learning can cost up to quadratic → for each point I need to compute all the distances

PROXIMITY BASED → outlier when NNs are far away (eg. 3-NN)

- effectiveness relies on the quality of the proximity measure → distance based or density based
- assumes that an outlier has a significantly higher proximity than the other data

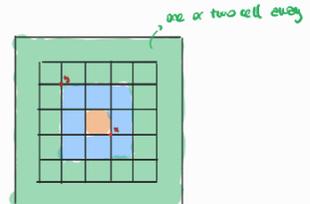
1) DISTANCE BASED: for each object $o \in D$ examine the # of objects in the r -neighbourhood ($r = \text{distance threshold}$)

• o is an outlier if $\frac{|\{o' \mid \text{dist}(o, o') \leq r\}|}{|D|} \ll r$ or k -th NN has $\text{dist}(o, o_k) > r$ with $k = \lceil \pi \cdot |D| \rceil$

Complexity is $O(n^2)$ but I can stop when I find $r \cdot |D|$ objects with distance $\leq r$ → linear average complexity

2) distance + grid based → data space is partitioned in hyper-cubes, each with a diagonal of $\frac{r}{2}$

- level 1: for any point x in cell C and y in a level 1 cell → $\text{dist}(x, y) \leq r$
- level 2: for any point x in C and y : $\text{dist}(x, y) \geq r \Rightarrow y$ is in a level 2 cell



not all made r -neighbourhood
 $b_1 = \text{total \# of objects in } L_1$
 $b_2 = \text{total \# of objects in } L_2$
 $a = \# \text{ of objects in } C$

- 1) pruning by L_1 : if $a + b_1 > \lceil \pi \cdot |D| \rceil \Rightarrow$ all objects in C are not outliers
- 2) pruning by L_2 : if $a + b_1 + b_2 < \lceil \pi \cdot |D| \rceil + 1 \Rightarrow$ all objects in C are outliers

The other objects I need to check manually

3) local distance based → useful to handle cluster with different densities → identify local outliers

$$D_k(x_i) = \frac{1}{k} \sum_{j \in N_k(x_i)} \|x_i - x_j\|$$

↑ avg dist k-NN
↓ k-NN of x_i

when $x \in \mathbb{R}^n$ with \uparrow very large $\|\cdot\|_2$ doesn't work well

$$O_k(x_i) = \frac{D_k(x_i)}{\frac{1}{k} \sum_{j \in N_k(x_i)} D_k(x_j)}$$

↑ outlieriness

$O_k(x_i) > 1 \rightarrow x_i$ is further away from neighbours than expected
 $O_k(x_i) \in [0, +\infty)$

- problem: when I have two close clusters with different densities (one sparse, the other dense), then the boundary of the dense one is considered outliers
- two parameters are needed (k and threshold for O_k)
- time complexity $\in O(n^2 \cdot \# \text{ dimensions})$



4) Density based outlier detection → density around an outlier is significantly different from the one around its neighbours

$$\text{dist}_k(o) = \text{distance of the } k\text{th NN}$$

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

more object with the same distance $\text{dist}_k(o)$

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

local reachability density

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

$k = \#$ of neighbours to take into account for density

local outlier factor

$$\text{LOF}_k(o) = \frac{\frac{\text{Lrd}_k(o)}{|N_k(o)|}}{\frac{1}{|N_k(o)|} \sum_{o' \in N_k(o)} \frac{\text{Lrd}_k(o')}{|N_k(o')|}} = \sum_{o' \in N_k(o)} \text{Lrd}_k(o') \cdot \frac{\text{reachdist}_k(o' \leftarrow o)}{|N_k(o)|}$$

→ compare local densities of the k-NN

↓ avg of the ratios $\text{LOF}_k(o) \leq 1 \rightarrow \text{OK}$

- same problem for the boundary as before

CLUSTERING BASED → normal data belongs to large and dense cluster, whereas outliers to small or sparse clusters or do not belong to any cluster
- clustering is expensive

1. x do not belong to any cluster → DBSCAN but two params are needed (ϵ , #points) → common to all the space (no local)

2. Far from its closest cluster → using k-means if $\frac{\text{dist}(o, c_o)}{\text{avg. dist}(c_o)} \gg 1 \rightarrow$ likely an outlier
↓ only spherical clusters

3. In small clusters (typically in large D, outliers makes cluster) eg. intrusion detection

b is the boundary between large and small clusters if one of the formula holds:
large cluster $LC = \{C_i \mid i \leq b\}$ and small as $SC = \{C_j \mid j > b\}$

$$\sum_{i \in LC} |C_i| \geq \alpha |D|$$

least large cluster

$$\frac{|C_b|}{|C_{b+1}|} \geq \beta$$

first small cluster

CBLOF (cluster based local outlier factor) is associated with each object

$$\text{CBLOF}(o) = \begin{cases} |C_i| \cdot \min\{\text{dist}(o, C_j)\} & \text{when } C_i \in SC: o \in C_i, C_j \in LC \\ |C_i| \cdot \text{dist}(o, C_i) & \text{when } C_i \in LC: o \in C_i \end{cases}$$

↑ small ↑ large

- first I need to cluster the data, then I need to find the CBLOF for $\forall o \in D \rightarrow$ high computational cost
- larger CBLOF(o) more likely it is that o is an outlier → small C_i for from LC or o_i far away from centroid

CLASSIFICATION BASED → need to face the high imbalance set (supervised methods)

1) ONE CLASS MODELS → learn the decision boundary of the normal class.

- can identify outliers that are not similar to those in the D (unseen anomaly)

a) Can be used to learn the decision boundary **SVDD** → construct an hypersphere around almost all positive class points with the minimum radius
- outlier if it falls outside the sphere

b) Using NN-Description (NN-d)

$$f_{\text{NN-d}}(z) = I\left(\frac{\|z - \text{NN}^{\text{tr}}(z)\|}{\|\text{NN}^{\text{tr}}(z) - \text{NN}^{\text{tr}}(\text{NN}^{\text{tr}}(z))\|}\right)$$

↑ test obj ↑ NN of z on training

$$I(x) = \begin{cases} 1 & \text{when } x \leq 1 \rightarrow \text{accepted} \\ 0 & \text{when } x > 1 \rightarrow \text{outlier} \end{cases}$$

$x = \frac{d_1}{d_2} \geq \frac{d_1}{d_2} > 1$

2) SEMI-SUPERVISED LEARNING → combining classification and clustering

- find a large cluster C and a small C_1 - since some obj in C carry the label "normal" → C is normal

- use a one-class model to identify outliers from C - since some obj in C_1 carry the label "outlier" → C_1 out
→ obj that doesn't belong to C or C_1

→ quality is based on the quality and representativity of the training set

- MINING CONTEXTUAL OUTLIERS

if the context is clearly identified I can find outliers in the limited dataset using conventional techniques.

If $D_{context}$ is not representative, then we can assume that similar context has similar behaviours (similar ages or locations)

Learn a mixture model on $U \rightarrow$ contextual variables and on $V \rightarrow$ behavioural attributes
 U has clusters U_1, \dots, U_n (on context str) V has clusters V_1, \dots, V_m

outlier score (based on context) \rightarrow if $0 \in U_j$ what p that $0 \in V_i$

$$S(o) = \sum_{U_j} p(o \in U_j) \sum_{V_i} p(o \in V_i) p(V_i | U_j)$$

I can model the normal behaviour based on the context \rightarrow I got an outlier when the prediction and the behaviour differs significantly
 - useful when I cannot identify clearly the context since I do not need to specify, but the model understands it autonomously \rightarrow Markov Chains
 eg. # of articles browsed in the part of one customer to be considered to understand if a new purchase is an outlier



- MINING COLLECTIVE OUTLIERS \rightarrow group of objects that deviates significantly from the entire dataset (more dense)

Individual points may seem normal and not outliers, but this change when I consider the entire group together
 \rightarrow time series may have a strange shape in an interval, but I couldn't understand it by a single point

I need to check the relations between data objects \rightarrow more difficult than other outliers (more computational complexity)
 bigger units than individual object/points

depends on the type of data:

- 1) temporal data needs to check smaller structured in time
- 2) spatial data explore local areas
- 3) graph data explore subgraphs

In contextual outlier analysis the structures to be analyzed are the context, here is not specified (needs to be learned / modeled)

Two methods:

- 1) Reducing the problem to classical outlier analysis \rightarrow identify structure units (time segments, local areas, subgraphs)
 Treat every structure units as a data object and extract features \rightarrow outlier analysis
- 2) Models the expected structural units directly \rightarrow eg. time series and Markov model, outlier when deviates
 - do not need to define structure units

- MINING OUTLIERS IN HIGH DIMENSIONS

problems:

- 1) Difficult to interpret since too many attributes (which subset particularly differs?)
- 2) data sparsity (curse of dimensionality) \rightarrow distance is dominated by noise

1) Extends conventional outlier detection
 HiOut Alg.

$w(o) = \sum_{i=1}^k \text{dist}(o, NN_i(o))$ \rightarrow top ϵ object with higher weight are outlier
 - I'm using ranks between distances and not absolute distances to make the outlier detection
 - Not good for interpretability

2) Dimensionality Reduction \rightarrow search outliers in sub-spaces (PCA)

GRID BASED \rightarrow project data on different subspaces and find areas where densities are much lower than the average

- 1) discretize data into grids with φ equi-depth (constant number of points $f = \frac{1}{\varphi}$) not equal width (fixed sized intervals)
 helps to cope with different distributions across attributes
- 2) search for cubes extremely sparse in k dimension (subspace)

If n objects in D , then the expected # of pts into a k -dimensional cube is $(\frac{1}{\varphi})^k n = f^k n$ \rightarrow Binomial random variable (or n or not)
 $\mu_k = n p$ $\sigma^2 = p(1-p)n$

The std of the # of pts in a k -cube is $\sigma = \sqrt{f^k(1-f^k)n}$
 sparsity coefficient $S(c) = \frac{n(c) - f^k n}{\sqrt{f^k(1-f^k)n}}$
 - If $S(c) < 0 \rightarrow$ less object than expected (more negative, more sparse)
 - assuming $S(c) \sim N$ we can get the significance level

\rightarrow Find the m cells with the lowest sparsity coefficient
 - brute force algorithm (check every cube in every subspace) $O(\varphi^k)$ for the \mathbb{R}^k subspaces (φ intervals per k dimension $\rightarrow \varphi^k$ cubes)
 - quality depends on grid resolution and position \rightarrow adopts a global approach on the expected # of points

3) Angle based methods \rightarrow angle measurements do not degrade with dimensions

$o \in D$ is an outlier if most objects are located in the same direction
 - use the distance-weighted angle variance as outlier score

ABOF(o) = $\text{VAR} \frac{\overline{ox} \cdot \overline{oy}}{\|o-x\| \cdot \|o-y\|}$ \rightarrow high ABOF means outlier Time complexity $O(n^3)$
 - angle based outlier factor



CLUSTERING HIGH DIMENSIONAL DATA → text documents or DNA arrays

Clustering should not only consider dimensions but also features (attributes)

- I can perform **feature transformation** to remove redundancies (PCA or SVD)
- or **feature selection** to find a subspace with relevant clusters

1) SUB-SPACE CLUSTERING → like in CLIQUE algorithm

clusters may exist only in some subspaces

They can be **similarity based** (bottom-up or top-down subspace search), **correlation based** (PCA) or BI-CLUSTERING

BI-CLUSTERING: cluster attributes and objects simultaneously

It has four requirements

- 1) only small set of objects participate in a cluster
- 2) A cluster only involve a small # of attributes] influenced by δ
- 3) objects can participate in more clusters or no cluster
- 4) attribute can participate in more clusters or no cluster

Example: customers and products purchased (attributes)

A bi-cluster is a sub-matrix with some consistent pattern. We have four types:

- 1) constant values $\forall i \in I, j \in J \Rightarrow e_{ij} = c$ (constant sub-matrix)
- 2) constant values on rows $\Rightarrow e_{ij} = c + a_i$
- 3) coherent values $\Rightarrow e_{ij} = c + a_i + b_j \Rightarrow e_{i_1 j_1} - e_{i_1 j_2} = e_{i_2 j_1} - e_{i_2 j_2}$ on columns but also on rows
- 4) coherent evolutions on rows $\Rightarrow (e_{i_1 j_1} - e_{i_1 j_2})(e_{i_2 j_1} - e_{i_2 j_2}) \geq 0$

Real world data is noisy. Two techniques to approximate bi-clusters:

a) optimization-based techniques → find one sub-matrix at the time with a greedy approach (local optimum)

δ -Bi-Clustering

mean of the i -th row $e_{i\cdot} = \frac{1}{|J|} \sum_{j \in J} e_{ij}$

mean of the j -th column $e_{\cdot j} = \frac{1}{|I|} \sum_{i \in I} e_{ij}$

all elements mean $e_{\cdot\cdot} = \frac{1}{|I||J|} \sum_{i \in I} \sum_{j \in J} e_{ij}$

The quality of a bi-clustering is:

mean squared residue

$$H(I \times J) = \frac{1}{|I||J|} \sum_{i,j} (e_{ij} - e_{i\cdot} - e_{\cdot j} + e_{\cdot\cdot})^2$$

→ from mean it deviates from a perfect coherent bi-cluster

if perfect $(c + a_i + b_j - c - a_i - b_j - c - \bar{a} - \bar{b}_j + c + a_i + \bar{b}_j)^2 = 0$

A sub-matrix $I \times J$ is a δ -bi-cluster if $H(I \times J) \leq \delta$ ($\delta = 0$ is a perfect bi-cluster with coherent values)

- A maximal δ -bi-cluster when \nexists another δ -bi-cluster that contains the first one

Since computation is costly, we use **greedy search** (approximate solution)

1) Deletion Phase → iteratively removes rows-columns while $H(I \times J) > \delta$

- to understand which row-column to remove at each step calculate the mean-squared-residue and remove the highest

$$d(i) = \frac{1}{|J|} \sum_{j \in J} (e_{ij} - e_{i\cdot} - e_{\cdot j} + e_{\cdot\cdot})^2 \quad \text{and} \quad d(j) = \frac{1}{|I|} \sum_{i \in I} (e_{ij} - e_{i\cdot} - e_{\cdot j} + e_{\cdot\cdot})^2$$

2) Addition Phase → extend the cluster until δ constraint is kept → add the smallest $d(i) / d(j)$

In this way I can find only one δ -bi-cluster (deterministic) → replace the elements of the cluster founded with random values (cannot be in more than one cluster)

b) Enumeration-based → use a tolerance threshold to specify the degree of noise allowed. Then try to enumerate all submatrices

δ -p-Cluster

for a 2×2 submatrix of $I \times J$

$$p\text{-score} \begin{pmatrix} e_{i_1 j_1} & e_{i_1 j_2} \\ e_{i_2 j_1} & e_{i_2 j_2} \end{pmatrix} = |(e_{i_1 j_1} - e_{i_2 j_1}) - (e_{i_1 j_2} - e_{i_2 j_2})|$$

if perfect $p\text{-score} = 0$

Allow a tolerance (noise) $p\text{-score} \leq \delta$ for every 2×2 sub-matrix of $I \times J$

p -score controls local noise in each point, while mean-squared residue is an average noise

Monotonicity property: if $I \times J$ is a δ -p-cluster → all the sub-matrices are δ -p-clusters

→ the similar techniques as FI analysis: for each combination J (features) find the maximal subset of objects I , such that $I \times J$ is a δ -p-cluster. → obj can be in more subsets

If $I \times J$ is not a sub-matrix of any found δ -p-cluster it is maximal

2) DIMENSIONALITY REDUCTION METHODS → sometimes is more efficient to use new features space rather than using original subspaces

SPECTRAL CLUSTERING

combine feature extraction and clustering → use the similarity matrix to perform dimensionality reduction

NJW Algorithm

1) calculate the affinity matrix $W_{ij} = e^{-\frac{\text{dist}(a_i, a_j)}{\sigma}}$ scaling parameter (how fast it decrease with the distance)

2) define a diagonal matrix $D_{ii} = \sum_{j=1}^n W_{ij} \rightarrow A = D^{-1/2} W D^{-1/2}$ → transformed sub-space $A = f(W)$

3) and k leading eigenvectors of A , project original data and cluster it with ordinary methods

↳ costly to compute

CONSTRAINED CLUSTER ANALYSIS → constrained like environment impediton that modify the distance between two points (a wall)

Constrained on instances:

- 1) must-link (x, y) → x and y on the same cluster
- 2) cannot-link (x, y) → on two different clusters eg. $\text{dist}(x, y) > d$

Constrained on clusters, like min # of objects, max diameter, shape (convex), # of clusters (k)

δ constrained (minimum separation)

$$\forall S_p, S_q \in S_i \times S_j, S_i \neq S_j \rightarrow \text{dist}(S_p, S_q) \geq \delta$$

→ equivalent to a must-link (x, y) : $\text{dist}(x, y) < \delta$

→ conjunction of constraints

ϵ constrained

$$\forall S_i: |S_i| > 1, \forall S_p \in S_i, \exists S_q \neq S_p, S_q \in S_i: \text{dist}(S_p, S_q) \geq \epsilon$$

→ constrained on the compactness

→ every point x must-link with a point y : $\text{dist}(x, y) \leq \epsilon$

→ disjunction of constraints

Constrained on similarity measures

- requirement that the similarity calculation must respect
- euclidean distance (for walking in a city is not good (walls))

Constraints can be **hard** (cannot be violated) or **soft** (better not to violate) → also called **PREFERENCES** → want to find clustering that maximally respects constraints

→ It needs to handle **conflicting** constraints (must-link and cannot-link of same x, y)

→ Measure the **usefulness** of a set of constraints:

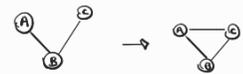
1) **Informativeness**: infos carried by the constrained beyond the clustering models

D (dataset), A (algorithm), C (constraints) → $I_A(C, D)$ = # of constraints violated during the clustering of D by A

2) **Coherence**: degree of redundancy among C

- helps to filter out useless constraints

- Adding constraints limits the search space of the algorithm → quicksize



1. **CoP-K-means** → handle hard constraints (given as part of the dataset labelled)

- for must-link compute the transitive-closure (add implied constraints). Then replace these objects with their mean → **Super-instance**

- modify the original K-means to compute the K means on the feasible objects (don't consider cannot-link)

2. Handle soft constraints: add a **penalty** when a constraint is violated

- then the objective will be to minimize an objective → sum of distances (like K-mean), adjusted by the penalties

• when must-link (x, y) but they are assigned different clusters, add $\text{dist}(c_1, c_2)$

• when cannot-link (x, y) but assigned on the same cluster, then add $\text{dist}(c, c')$
↳ cluster of x
 ↳ closest cluster to C that can accommodate x or y

Can be complex to find these clusters ⇒ optimization methods based on the domain of the application

• eg. group closed points in micro-clusters → reduce complexity instead of handling single points

GRAPH CLUSTERING

Network → a collection of entities that are interconnected with links (eg. social networks, communication networks, biological networks)

Graph → mathematical representation of a network (nodes and edges) $G(N, E)$
 - edges and origin: visual inspection doesn't work

To cluster graphs we can adopt two techniques:

- 1) Introduce a similarity metrics → standard cluster algorithms on dissimilarity matrix (high dimensional clustering if large graphs)
- 2) Algorithms that works directly on graphs

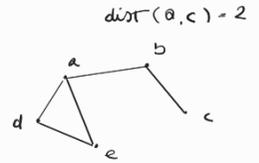
SIMILARITY METRICS

1) **Geodesics distance** $\text{dist}(A, B) = \#$ of edges of the shortest path between A and B

Eccentricity of a node $\text{eccen}(v) = \max_{u \in V - \{v\}} \text{dist}(v, u)$

Radius of graph $r = \min_{v \in V} \text{eccen}(v)$ → distance between the most central point and the furthest border

Diameter of graph $d = \max_{v \in V} \text{eccen}(v)$ → largest distance of any pair of vertices
 - A **peripheral vertex** is a vertex that achieves the diameter (d, e, c)



The geodesics distance doesn't measure well the closeness of vertices in particular domains

- Bob and Alice are vertices in a network, the shortest distance doesn't account for all the other possible interactions between the two

The similarity in a social network can have two meanings

1. **Structural context-based similarity**: two customers are similar if common neighbors
 - people that receive common suggestion from multiple friends often make the similar decision
2. **Random walk similarity**: a company sends promotional infos to both Ada and Bob. They randomly forward to their friends. The closeness is the likelihood that other customers simultaneously receive the promotion

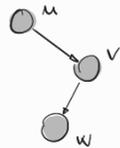
2) **SimRank**: context based similarity. Directed graph

$$I(v) = \{u \mid (u, v) \in E\}$$

↑ IN-NEIGHBOURHOOD

$$O(v) = \{w \mid (v, w) \in E\}$$

↓ OUT-NEIGHBOURHOOD



$$s(u, v) = \frac{c}{|I(u)||I(v)|} \sum_{x \in I(u)} \sum_{y \in I(v)} s(x, y)$$

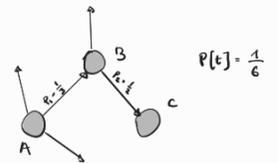
- If no neighbors $s(u, v) = 0$
- recursive formula

Compute it iteratively with a matrix $S_i(*, *) \in [0, 1]^{n \times n}$

$$s_0(u, v) = \begin{cases} 0 & \text{if } u \neq v \\ 1 & \text{if } u = v \end{cases} \rightarrow \text{A lower bound on the actual value}$$

$$\begin{cases} s_{i+1}(u, v) = \frac{c}{|I(u)||I(v)|} \sum_{x \in I(u)} \sum_{y \in I(v)} s_i(x, y) & \text{if } u \neq v \\ s_{i+1}(u, v) = 1 & \text{if } u = v \end{cases} \rightarrow \text{non decreasing}$$

Fix a number of iteration to get the approximation $O(K n^2 d_2)$



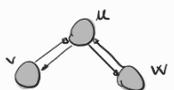
3) **Random Walk**

$$d(u, v) = \sum_{t: u \rightsquigarrow v} P(t) \ell(t)$$

$$P(t) = \begin{cases} \frac{1}{\prod_{i=1}^{t-1} |O(u_i)|} & \text{if } \ell(t) > 0 \\ 0 & \text{if } \ell(t) = 0 \end{cases}$$

$d(u, v)$ is the number of steps a random surfer needs to take to go from u to v
 - not symmetric

Expected Meeting Distance (EMD) is the expected # of steps two surfers needs to take before meeting each other - symmetric



$$\begin{aligned} m(u, v) - m(u, w) &= +\infty \\ m(v, w) &= 1 \end{aligned}$$

To formally define EMD we use the derived G^2 of node-pairs

- each node $(u, v) \in V^2$ can be thought as the present state of two random surfers in G
- an edge in E^2 $(u, v), (t, w)$ means that $u \rightarrow t$ and $v \rightarrow w$, possible steps of two surfers
- A tax in G^2 having length n also in G has length n

EMD is the expected distance in G^2 from (u, v) to any node $(x, x) \in V^2$ (meeting point)

- If G is strongly connected, G^2 may not be → no path from $(u, v) \rightsquigarrow (x, x) \in V^2$, then $m(u, v) = +\infty$
- If from (u, v) there is another path to (u, v) there is a problem on the definition if passes through a meeting point $(x, x) \in V^2$

$$m(u, v) = \sum_{t: (u, v) \rightsquigarrow (x, x)} P(t) \ell(t)$$

EXPECTED f-MEETING distance: Map all distances to a finite interval, instead of computing $E(t)$, we can compute $f(E(t))$
 $f: [0, \infty) \rightarrow [a, b]$ (non-negative, monotonic, bounded) - solution to infinite EMD (when G^t is not strongly connected)

$$s'(u, v) = \sum_{t: (u, v) \rightarrow (x, x)} P(t) c^{E(t)}$$

$$c \in (0, 1)$$

$$s'(a, b) = 0 \rightarrow \text{no path between } a \text{ and } b$$

$$s'(a, b) = 1 \rightarrow a = b$$

Equivalence to simrank with parameter c

$$s'(a, b) = \sum_{x=1}^{10^{(a,b)}} \sum_{t: 0, (a,b) \rightarrow (x, x)} P(t) c^{E(t)+x} = \frac{c}{|0(a)| |0(b)|} \sum_{i=1}^{10^{(a,b)}} \sum_{j=1}^{10^{(a,b)}} s'(0_i(a), 0_j(b))$$

It's a random similarity but also a structural context-based similarity (also sim-rank)

GRAPH CUT → make a cluster points should be well connected, instead in different clusters connected in a much weaker way

A cut $C(S, T)$ is a partitioning, where $V = S \cup T \wedge S \cap T = \emptyset$

The cut set is $\{(u, v) \in E \mid u \in S, v \in T\}$

Size of the cut is the number of edges in the cut set, or if it is a weighted graph, the sum of weights

1) **MINIMUM CUT** → $C(S, T)$, where the cut size of T is not greater than any other cut's size

- it takes polynomial time to compute

- is not that good because it will find small clusters (even a single point)

2) **SPARSEST CUT** → Every node $n \in$ cut set has most of the edges belonging to one cluster

Given a cut $C(T, S)$

$$\phi = \frac{\text{cut size}}{\min(|S|, |T|)}$$

→ the sparsest is the cut with the minimum ϕ

favours balanced (bisection) and sparse clusters
 Thanks to min

- the problem is NP hard but $O(\sqrt{E})$ approximation

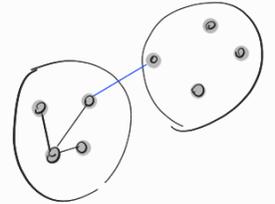
Theoretically graph clustering problems can be regarded as optimal cut problems but they are computationally expensive

MODULARITY of CLUSTERING

$$Q = \sum_{i=1}^n \left[\frac{E_i}{|E|} \left(\frac{d_i}{2|E|} \right)^2 \right]$$

of edges in the i-th cluster → completely inside C_i
 sum of degrees of i-th cluster vertices
 probability random edge falls in C_i
 probability e entirely in C_i

- High Q means more edges in clusters than you expect by chance (randomly cutting the graph)
- Optimal cut has the highest modularity



SCAN ALGORITHM

define $\Gamma(v)$ as the immediate neighbourhood of a vertex $v \in V$

$$\Gamma(v) = \{u \mid (v, u) \in E\} \cup \{u\}$$

$$O(v, w) = \frac{|\Gamma(v) \cap \Gamma(w)|}{|\Gamma(v)| + |\Gamma(w)|} \rightarrow \text{structural similarity: high for members of a clique, small for hubs and outliers}$$

ϵ -Neighbour of v : $N_\epsilon(v) = \{w \in \Gamma(v) \mid O(w, v) \geq \epsilon\}$

v is a core vertex $\Leftrightarrow |N_\epsilon(v)| \geq \mu$

problem: select the two variables

$$\text{Dir RECH}_{\epsilon, \mu}(v, w) \Leftrightarrow \text{CORE}_{\epsilon, \mu}(v) \wedge w \in N_\epsilon(v)$$

If a vertex $v \in N_\epsilon(u)$ and u is a core vertex → v is assigned to the same cluster of u

growing process until no cluster can be enlarged

Structure Reachable

→ transitive closure of direct reachability. A vertex v can be structure reachable from $u \Leftrightarrow \exists \langle w_1, \dots, w_k \rangle$, w_1 reachable from u , w_i reachable from w_{i-1} , v reachable from w_k . $\text{RECH}(u, v)$

Structure connected

→ v, w are connected if $\exists u$ core object: v and w can be reached from u
 $\text{CONNECT}_{\epsilon, \mu}(v, w) \Leftrightarrow \exists u \in V: \text{RECH}_{\epsilon, \mu}(u, v) \wedge \text{RECH}_{\epsilon, \mu}(u, w)$

A vertex is an **hub** if in $\Gamma(v)$ the vertices are in different clusters and v doesn't belong to any cluster

All other point that are not in any cluster and are not hubs are **outliers**

- **Clique** is a subset of V where each v_i, v_j are adjacent → tight social group (over used in slightly relaxed way)

A cluster C is **connected** if $\forall v, w \in C \rightarrow \text{CONNECT}_{\epsilon, \mu}(v, w)$

A cluster C is **maximal** if $\forall v, w \in V: v \in C \wedge \text{RECH}_{\epsilon, \mu}(v, w) \Rightarrow w \in C$

SEQUENTIAL PATTERN MINING → discovering of subsequences in a sequence DB

format [transactionID, <Ordered Sequence of Events>] eg. [T1, <(bread, milk), (bread, milk, sugar)>] → order between tuples of the sequence
 [T2, <(bread), (milk)>]

Web usage mining → each event is a single item (1-itemset) since a user can access one page at the time

- I want to discover users' navigational patterns
- 1. server-side: multiple users - single server. Good for recommendation systems (eg. Amazon). The cache is a problem
- 2. Client-side: single user - multiple servers. Web content personalization. Solves cache problems
- 3. Proxy-server: multiple users - multiple servers

Sequential DB is a set of sequential records (ordered sequences).

We want to find the set of all frequent sequences S in D, of items I at the given min-sup

A **sequence** $S = \langle e_1, e_2, \dots, e_n \rangle$ with e_i an itemset
 e_j is a lexicographically ordered list of items

eg. (abc) < (abcd) and (ab) < (aba)

$$t = \{i_1, \dots, i_k\}, t' = \{j_1, \dots, j_l\} \rightarrow t < t' \iff (\exists 0 \leq k < \min\{k, l\} : i_k = j_k \ \forall j < k \wedge i_k < j_k) \vee (k < l \wedge i_k = j_k \ \forall j < k)$$

TOTAL ORDER for lexicographically ordered itemsets

A item can appear only once inside an itemset, but more times across itemsets of a sequence

- A sequence $\alpha = \langle e_1, \dots, e_m \rangle$ is a **subsequence** of another sequence $\beta = \langle e_1, \dots, e_n \rangle$ ($\alpha \leq \beta$) if $i_1 < \dots < i_m$ and $i_j \geq 1 \wedge i_m \leq n : e_{i_j} \leq e_{i_j}$
- $\langle ad \rangle \leq \langle a(bd)bcd \rangle$
- a sequence is frequent if it is a subsequence of at least min-sup sequences in D
- a subsequence is **maximal** if it's not a subsequence of any other sequence
- S is frequent **closed** sequence if $\nexists S_0 : S_0 \supset S \wedge \sigma(S_0) = \sigma(S)$

$$\sigma(S) = \frac{|S' : S \supset S'|}{|D|}$$

support

ALGORITHM

- Sort Phase** → converts the original transaction DB to a sequential DB
- Itemset Phase** (large itemset) → find only the large itemsets (Apriori property, otherwise the sequence containing that itemset cannot be frequent)
 - has more than min-sup occurrences, counting at most one per customer (different from traditional methods)

Then **map** each large itemset to a set of contiguous integers
 - constant time compare of two itemsets and reduce the time to check if $S_a \leq S_b$

- Transformation Phase**
 - replace each transaction with all the itemsets contained in the transaction
 - empty transactions (with no itemset) are dropped, but considered for support counting?

ITEMSET	MAP
{30}	1
{40}	2
{70}	3
{40 70}	4
{90}	5

CUSTID	ORIGINAL SEQ	TRANSFORMED	MAPPED
1	<(10 20) (30) (40 60 70)>	<{30}, {40}, {70}, {40 70}>	<{1}, {2, 3, 4}>

(10 20) → dropped
 (40 60 70) → {(40), {70}, {40 70}} → {2, 3, 4}

- Sequence Phase** → Similar approach as Apriori
 Different algorithms:
 - Count-all: count all large sequences, including non-maximal (Apriori-All)
 - Count-some: try to avoid counting non-maximal S (Apriori-Some)

- Maximal Phase** → find maximal sequences among large sequences

APRIORI ALL 1) $\forall S_1, S_2 \in L_{k-1} : S_1, S_2$ have in common the first $[1, k-2]$ elements, join them (add the last element of S_2 to S_1)

$$L_3 = \{ \langle 123 \rangle, \langle 234 \rangle, \langle 124 \rangle, \langle 134 \rangle, \langle 135 \rangle \}$$

$$L_4 = \{ \langle 1234 \rangle, \langle 1243 \rangle, \langle 1345 \rangle, \langle 1345 \rangle \}$$

Important to add both of them, since I could join the two subsequences in both ways (the order is important for a sequence)

2) Delete sequence $S \in L_k$ if $\exists S' : S' \supset S$ (Apriori property)
 $C_4 = \{ \langle 1234 \rangle, \langle 1243 \rangle, \langle 1345 \rangle, \langle 1345 \rangle \} \Rightarrow C_4' = \{ \langle 1234 \rangle \}$

APRIORI SOME → two phases 1) **Forward Phase**: Find all large sequences of some lengths, determined by the **next()** function

eg. Apriori-All next(k) = k+1
 - skip some lengths if I think the corresponding sequences are frequent and non-maximal
 (I could use $next_k = \frac{|L_k|}{|C_k|}$ as an index)

- In candidate generation, when L_{k-1} is not available I will use C_{k-1} to generate C_k
 - not a problem in terms of correctness, since $L_{k-1} \subseteq C_{k-1}$

2) **Backward phase**: for all lengths that we skipped we delete the sequences $\in C_k'$ that are contained in some large sequence

- count the remaining sequences and keep the large
- Also delete, from the sequences found in the first phase, the one that are not maximal

APRIORI DYNAMIC SOME

1) **initialization phase** → All $S: |S| \leq \text{step}$ are counted (Here is not only next(k) as in Apriori Some)

2) **forward phase** → all sequences whose lengths are multiples of steps are counted
 - We can directly find the C_7 sequences by joining $L_2 \bowtie L_3$ in a different way ← ?

●: classic Apriori generation $C_k = L_{k-1} \bowtie L_{k-1}$

●: on the fly generation (OTF) $C_{k+step} = L_k \bowtie L_{k+step}$

3) **Intermediate phase** → generate the skipped candidate set

for each step 4) **Backward phase** → count the candidate set and remove the non maximal

Example step=3 → assuming L_1 is found empty I have $C_1, L_1 = C_2, L_2 = C_3, L_3 = C_4, L_4 = \{ \}$

3a) Generating C_7 from L_6 and C_8 from C_7

3b) Count C_8 and find L_8 , remove non maximal from C_7 , then count and find L_7

Repeat for C_4 and C_5

OTF generation → takes as argument L_k, L_{step} and a customer sequence c and returns the candidate $(k+step)$ -sequences $\leftarrow c$

$c = \langle c_1 \dots c_n \rangle$ If $s_k \in L_k$ and $s_{step} \in L_{step}$ are contained in c , but do not overlap → possible candidate subsequence of len $k+step$

$$s_k \text{ end} = \min \{ j : s_k \leq \langle c_1 \dots c_j \rangle \}$$

$$s_{step} \text{ start} = \max \{ j : s_{step} \leq \langle c_j \dots c_n \rangle \}$$

$$\Rightarrow \forall s_k, s_{step} \quad \text{join } s_k \bowtie s_{step} \iff s_k \text{ end} < s_{step} \text{ start}$$

→ It could be better to find C_{k+step} from L_k and L_{step} , rather than modify the apriori and find it by $L_k \bowtie L_k$ (with first k -step equal elements)

- Anyway the performance test shows that Apriori Dynamic Some is the **slowest** because it generates too many candidates - there is no non-maximal pruning in the forward phase

APRIORI ALTERNATIVES Bottlenecks of Apriori Methods

1) Many DB scan

2) exponential # of short candidates when mining long pattern: a 100-sequential pattern requires $\sum_{i=1}^{100} \binom{100}{i} \approx 10^{10}$

f-list: list of frequent items (i-itemset) in support descending order

frequent f-list = b:5, c:4, a:3, d:3, e:3, f:2 ← sorted L_1 , 1 DB scan to find it

PROJECTION RULE FOR SEQUENTIAL PATTERNS

i-projected db → set of sequential patterns containing i , but no items that followed i in the f-list
 ↳ having i as subsequence

S is projected in the i -projected db as s_i , by removing infrequent items and any item following i in f-list

$$\langle (ab)(bf)abf \rangle \xrightarrow{b\text{-proj}} \langle bb \rangle$$

DB PARTITION I can divide the dataset in 6 **non-overlapping** subsets

pattern sequences containing f
 pattern sequences containing e , but not f
 pattern sequences containing d , but not f and e
 ⋮

Two way of projecting the db in order to achieve something more manageable (RAM resident) → final result is the **same**:

1) **parallel projection**: generate all the projected db at a time but its result is a set **larger** than original db

- for each sequence, project it to all the elements of the f-list, following the rule → added to more db
- the transaction can be projected in parallel in every db
- Imagine to have on avg l frequent item per transaction → projected into $l-1$ db's (excluding the first)
- The total size is $1 + \dots + (l-1) = \frac{l(l-1)}{2}$ → size (DBs) = $\frac{l(l-1)}{2}$ size (DB)

2) **Partition Projection** project a sequence to the projected db of the **least** frequent item in it, following the rule

↳ because a sequence is added to only one DB → partition the sequences (not / but serial)

- After the processing of the db, remove the least frequent element and continue the projection
- Size of the projection always smaller than the original db - still one DB scan

projection level \rightarrow size of the subsequence the sequence in DB must have

Two ways of mining projected db

- 1) Level-by-Level \rightarrow find f.-ext with one scan (L1), partition the db into f.-ext, mine L2, partition the projected dbs to L2,
 - 2) Alternative level \rightarrow postpone the generation of projected db using annotations, free-span algorithm
- } both recursive

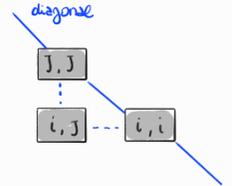
FREE SPAN \rightarrow similar to FP-growth \rightarrow maximum 3 scans of the original db

b	4					
c	(4,3,0)	1				
a	(3,2,0)	(2,1,1)	$\rightarrow \langle a, a \rangle$			
d	(2,2,2)	(2,2,0)	(1,2,1)	1		
e	(3,1,1)	(1,1,2)	(1,0,1)	(1,1,1)	1	
f	(2,2,2)	(1,1,0)	(1,1,0)	(0,0,0)	(1,1,0)	2
	b	c	a	d	e	f

- 1) first scan get the f.-ext
- 2) construct the **Frequent Item Matrix** with another DB scan \rightarrow triangular matrix $|L1| \times |L1|$
 - $F[i,j] \rightarrow$ has one counter, reading the appearance of $\langle i, j \rangle$ in the db
 - $F[j,k] = (A, B, C)$ with $A: \langle i, i \rangle, B: \langle i, j \rangle, C: \langle i, k \rangle$

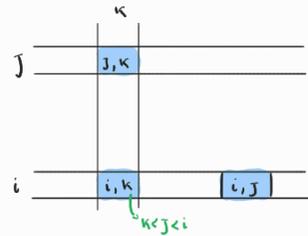
With this matrix we can:

- generate L2 \rightarrow If a counter is $\geq \text{min_sup} \Rightarrow$ add the corresponding sequence to L2
 - generate annotations for item repeating patterns
 - if $F[j,j] > \text{min_sup} \rightarrow$ generate $\langle j j^* \rangle$
 - for column $j \neq i$:
 - if $F[i,i] \geq \text{min_sup} \rightarrow i^*$
 - if $F[j,j] \geq \text{min_sup} \rightarrow j^*$
 - if only one of the tree counters of $F[i,j]$ is above min-sup $\rightarrow \langle i^* j^* \rangle$ or $\langle j^* i^* \rangle$ otherwise generate $\langle i^* j^* \rangle$
- c) generate annotations for projected dbs of the form $\{a_i a_j\} \cdot \{b_1, \dots, b_q\}$



take a triangle with one vertex (i, j) $j < i$, then choose $k < j$. If $F(i, j), F(i, k), F(j, k)$ forms a **generating triple** add $\{i, j\} : \{k\} \rightarrow$ compute for all k and make the union

$F(e, f) = (1, 1, 0)$ $F(a, e) = (3, 0, 1)$ $F(a, e) = (1, 1, 3)$
 \rightarrow it's not a generating sequence, because $\langle efa \rangle, \langle aef \rangle, \langle eaf \rangle, \dots$ cannot be frequent
 But if $F(a, e) = (3, 1, 1)$ it could generate the pattern $\langle afe \rangle$
 \rightarrow adding $\langle fe \rangle \cdot \{a\}$



- 3) Discard F matrix, then scan the db the third and last time in order to find item-repeating sequences and projected dbs
- 4) Further mine the projected dbs \rightarrow If they contains exactly tree items, then with one one scan we can generate all otherwise recursive, start from one

TIME SERIES → set of data points ordered in time

They can be decomposed in :
 · trend : slow moving change
 · seasonality
 · residual : whatever cannot be explained by trend or seasonality

Forecasting → predicting the future using historical data

· It's different from regression tasks because here points have an order that could hide some pattern like seasonality

Mean Absolute Percentage Error :
$$MAPE = \frac{1}{n} \sum_{i=1}^n \left| \frac{A_i - F_i}{A_i} \right| \cdot 100$$

A **baseline model** (trivial solution to the forecasting problem, in order to make some benchmark) could be :

- 1) use the last known value → no mean or seasonality taken into account **NAIVE METHOD**
- 2) use the time serie mean
- 3) naive seasonal forecast → takes the last observed cycle and repeat it

Random Walk Process : process in which there is an equal probability of going up or down by a random number (ϵ_t)

$y_t = C + y_{t-1} + \epsilon_t$ with $\epsilon_t \sim N(0, \sigma^2)$ also called white noise

· Formally a Random walk is a time serie whose first difference ($C + \epsilon_t$) is stationary and uncorrelated

$E[y_t] = E[C] + E[y_{t-1}] + E[\epsilon_t] = tC + y_0$
 $Var[y_t] = t\sigma^2$

To forecast random walk on the long period we cannot use statistical learning, since they change randomly

I can use naive methods or baselines **Drift method** :
$$\hat{y}_{t+k} = \frac{y_t - y_0}{t} \cdot k$$

Stationary time serie : $\{y_t\} \rightarrow \forall t, t'$ the distribution of (y_t, \dots, y_{t+k}) does not change

- constant mean
- constant variance
- constant auto-correlation $cov(y_t, y_{t+k}) = \delta(k) \forall t$ → could be correlated, but independently from time

· To make a time serie stationary we can implement some **transformation** technique :

- differentiating $y'_t = y_t - y_{t-1}$ helps stabilizing the mean
- applying the log helps stabilizing the variance → e^x to untransform

AUGMENTED Dickey-FULLER TEST → assume the time serie is well described by $y_t = C + \alpha_1 y_{t-1} + \epsilon_t$ **AR(1)**

α_1 is called the **root** of the time serie → $\{y_t\}$ will be stationary ↔ $\alpha_1 \in (-1, 1)$ → **NO UNIT ROOT** (cannot be negative)

· simple to see since $Var(y_t) = \alpha_1^2 Var(y_{t-1}) + \sigma_\epsilon^2 \Rightarrow \sigma^2 = \alpha_1^2 \sigma^2 + \sigma_\epsilon^2 \Rightarrow \sigma^2 = \frac{\sigma_\epsilon^2}{1 - \alpha_1^2}$

· **NH** (null hyp) : a unit root is present (not stationary) → check the returned p-value with significance level
 · p-value < α_{signif} → reject

Auto correlation $r_k = \frac{1}{T} \sum_{t=1}^{T-k} (y_t - \bar{y})(y_{t+k} - \bar{y})$ $r_k = \frac{C_k}{C_0} = \frac{cov(y_t, y_{t+k})}{Var(y_t)}$

AutoCorrelation Function (ACF) $f(k) = r_k$
 · when we have trend, autocorrelation for small k tends to be large and positive
 · with seasonality, we have peaks at multiples of the seasonal frequency

→ stationarity ⇒ white noise wn is a special case with zero autocorrelation $E[\epsilon_t \epsilon_{t-1}] = 0$

→ stationary, but there is autocorrelation, plot ACF coefficients and if they are not significant after lag q ⇒ MA(q)

MOVING AVERAGE MODEL $MA(q) = \mu + \epsilon_t + \theta_1 \epsilon_{t-1} + \dots + \theta_q \epsilon_{t-q}$
 · depends linearly with mean, current and past error terms
 · it is a stationary time serie

- to choose q, I need to plot ACF and see after which lag q, all coefficients are not significant
 · predicting after q steps in time simply returns the mean (no error to use) → use **rolling forecast** to predict up to q steps ahead, then retain with new observation

when ACF is not conclusive, then plot PACF and if coefficients are abruptly non significant after lag $p \Rightarrow AR(p)$

AUTOREGRESSIVE MODEL

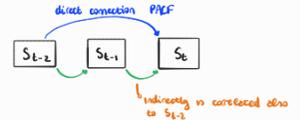
\rightarrow regression of a variable against itself

$\cdot C$ is no longer the mean

$$AR(p): y_t = C + \phi_1 y_{t-1} + \dots + \phi_p y_{t-p} + \epsilon_t$$

\cdot to find the order p I need to calculate the **partial auto correlation function (PACF)**

- PACF measures the relation after removing the effects of the observation n between
- is actually a regression to find ϕ_p



AUTOREGRESSIVE MOVING AVG

\rightarrow when both ACF and PACF doesn't rapidly decay

$$ARMA(p, q): y_t = C + \phi_1 y_{t-1} + \dots + \phi_p y_{t-p} + \epsilon_t + \theta_1 \epsilon_{t-1} + \dots + \theta_q \epsilon_{t-q}$$

$\cdot ARMA(0, q)$ is equivalent to $MA(q)$, $ARMA(p, 0) = AR(p)$

\cdot To determine p and q use the **Analytic Information Criterion (AIC)**

$$AIC(p, q) = 2k - 2 \ln(\hat{L}) = p + q - 2 \ln P(y_t | ARMA(p, q))$$

↑ explainability ↑ goodness on fitting data

1) for every tuple (p, q) evaluate AIC and select the best (the lowest)

2) Analyze the model residue which is the difference of $\Delta y_t = |y_t - \hat{y}_t|$

If they are **white noise**:

- \cdot uncorrelated \rightarrow Ljung-Box test
- $\cdot N(0, \sigma^2) \rightarrow Q-Q$ plot is a straight line

\Rightarrow use them to forecast, otherwise change (p, q)

Ljung-Box test

\rightarrow statistical test to determine if the autocorrelation is significant different from zero

H_0 : data is independently distributed (no auto-correlation)

p -value $> 0.05 \rightarrow$ cannot reject \rightarrow white noise

AUTOREGRESSIVE INTEGRATED MOVING AVG

\rightarrow when the serie is not stationary I will apply d times the differentiation and then ARMA

Excess # to make it stationary

$$y'_t = C + \phi_1 y'_{t-1} + \dots + \phi_p y'_{t-p} + \epsilon_t + \theta_1 \epsilon'_{t-1} + \dots + \theta_q \epsilon'_{t-q}$$

↑ differenced time serie ↑ differenced error

SEASONAL ARIMA

\rightarrow adds seasonality parameters to $ARIMA(p, d, q)$

$\rightarrow SARIMA(p, d, q)(P, D, Q)_m$

seasonal frequency = # of observation per cycle

P : order of seasonal $AR(p)$ process, D : order of seasonal integration, Q : order of seasonal $MA(Q)$ process

M-12 if 1 year cycle, with data every 4 month

$SARIMA(p, d, q)(0, 0, 0)_m = ARIMA(p, d, q)$

P-2 including two past value, multiple of $m \rightarrow y_{t-2}, y_{t-2m}$

- \cdot Now I need to fit (p, q, P, Q) tuples
- \cdot d and D are selected before

SARIMAX MODELS

\rightarrow it adds a linear combination of exogenous variable to SARIMA

$$y_t = SARIMA(p, d, q)(P, D, Q)_m + \sum \beta_i x_t^i \rightarrow \text{need to know the variables at time } t$$

\cdot to predict at $t' > t+1$ we need to forecast also the exogenous variables (may be easier to predict)

ML MODELS for time series

- \cdot statistical approaches works well when we have small dataset (less than 10k points) otherwise they are very slow
- \cdot ML when we have large datasets or residues are not white noise

- 1) Single-step model: predicts the single value for the next time step \rightarrow single step
- 2) Multi-step model: predicts the single value but for many time steps into the future
- 3) Multi-output model: when predicting for more than one variable (temp and wind speed) \rightarrow multiple single steps

THE ENCODING \rightarrow since we want to encode info for cyclical behaviour (if it increase always is not useful)

If I want to map the 24h cycle: $\vec{E} = (t_1, t_2)$

$$t_1 = \sin\left(t - s \frac{2\pi}{24 \cdot 60 \cdot 60}\right) \quad t_2 = \cos\left(t - s \frac{2\pi}{24 \cdot 60 \cdot 60}\right)$$

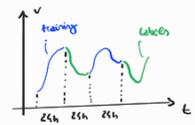
↑ timestamp in seconds

not sufficient 12am = 12pm

Data windowing process → we define a sequence of points that are training and those that are labels

ex. 24h period to predict the next 24h period

• huge data waste

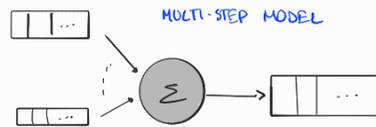
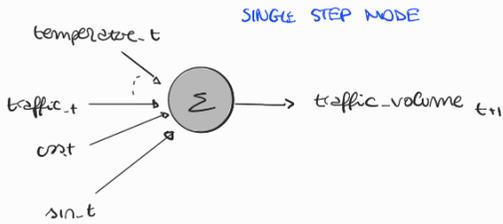


- **Sliding windows with step k**: shift of k positions the training/labels windows

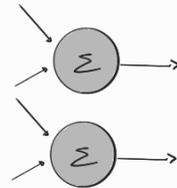
Baseline model for multistep process: input width = N, label width = N, $shift(k) = N$

1) last known value for the next N steps

2) repeat the last N steps to the next



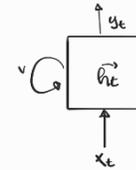
MULTIVARIATE MODEL



• simple multivariate linear regression

Often **RNN** are used → but they suffer of short-term memory

• previous values stop having influence in the far future (seasonality)



LSTM try to solve this problem introducing input and forget gate to select the new infos to store and the one to forget

FOUNDATIONAL MODELS → big pre-trained models

• normally good performances if they are trained with something similar

• nevertheless they are not explainable and may not be the most efficient solution

DATA STREAMING

Stream \rightarrow potentially unbounded, ordered sequence of instances

$$S = \{x_1, \dots, x_N\}$$

$$x_i \in X^d$$

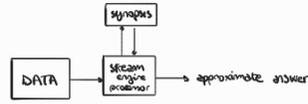
$$N \rightarrow +\infty$$

- traditional techniques \rightarrow entire dataset that needs to be scanned multiple times and randomly accessed. Computationally heavy learning phase
- data streams techniques \rightarrow impractical/impossible to store the entire db, multiple scans or random access

Typically we need to handle massive volumes of data and their rapid production frequency (low complexity training, impractical for each instance to re-execute the algorithm)

Requirements

- 1) single pass: record examined at most once
- 2) bounded storage: limited memory for storing synopsis (summarization of data)
- 3) real-time



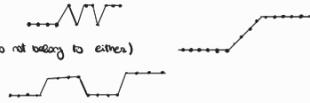
- 1) approximate answers
 - 2) a.a. with deterministic bounds $(\epsilon - \delta)$
 - 3) a.a. with probabilistic bounds $(\epsilon - \delta)$
- with $1 - \delta$ probability
distance from correct ans

SYNOPSIS \rightarrow particular data structure enable to incrementally summarize data

- 1) feature vectors \rightarrow summary of data instances by features extraction
- 2) prototype arrays \rightarrow keep only representative instances
- 3) coresets trees \rightarrow keep the summary in a tree structure (like BIRCH clustering algorithm)
- 4) grids \rightarrow keeps data feature space in a quantized grid

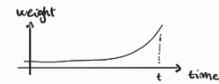
CONCEPT DRAFTS: unforeseen change in statistical properties of data stream instances over time

- 1) sudden concept drift: sudden change in two consecutive instances at once, after this time only instances of the new class are received
- 2) gradual concept drift: # of one class decrease gradually, the other increasing
- 3) incremental concept drift: evolving to the new class gradually (in the meantime do not belong to either)
- 4) recurring concept drift: data instances change between two classes in turns



WINDOW MODEL \rightarrow helps in reducing concept drift effects

- 1) hamped window model: recent have more weight than older one (importance decreases over time). decay function to model the feature importance like: $f(t) = e^{-\lambda(t_{now} - t)}$
- 2) Land-mark window model: whole data in a window between two landmarks with same importance. consecutive windows do not intersect
- 3) Sliding window model: window shift one position at each time (FIFO style). consecutive windows mostly overlap

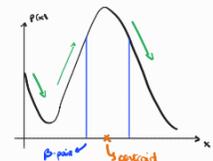


x_i is in the $m \cdot \lfloor \frac{t}{w} \rfloor$ window

CLUSTERING \rightarrow since no y_{true} is needed is easily used with data streams

ADAPTIVE STREAMING K-MEANS

- 1) collect e data instances
- 2) keep best clustering parameters (k , centroids) by evaluating silhouette coefficients $e O(d \cdot e \cdot k \cdot cs)$
 - determine the distributions (pdf) of each variable independently \rightarrow may follow different distributions
 - to find k_i identify directional change in the pdf: $e O(e)$ $k \in [k_{min}, k_{min} + k_{max}]$ for all features $e O(d \cdot e)$
 - split the pdf in k_i equiprobable areas (boundaries are called β -points), assign candidate centroids to the center of those regions (typically high data density)
- 3) continuous clustering with k means of the data stream, if a concept drift is detected, the restart the initialization phase (assign α single point) $e O(d \cdot k)$
 - \hookrightarrow predicted by comparing μ and σ stored



for each point e , calculate all distances $e O(d \cdot k)$ and repeat for cs different centroids sets

The worst case scenario to cluster a point, involving the computation of initialization phase is $e O(d \cdot e) + O(d \cdot e \cdot k \cdot cs) + O(d \cdot k)$