$$support(X \rightarrow Y) = \frac{count(X \cup Y)}{n}$$

$$confidence(X \rightarrow Y) = \frac{count(X \cup Y)}{count(X)}$$

# 1   Apriori frequent itemset algorithm

init-pass obtains the items in lexicographic order.

```
Algorithm Apriori(T)
    C₁ ← init-pass(T);
    F₁ ← {f | f ∈ C₁, f.count/n ≥ minsup};    // n: no. of transactions in T
    for (k = 2; F_{k-1} ≠ ∅; k++) do
        C_k ← candidate-gen(F_{k-1});
        for each transaction t ∈ T do
            for each candidate c ∈ C_k do
                if c is contained in t then
                    c.count++;
            end
        end
        F_k ← {c ∈ C_k | c.count/n ≥ minsup}
    end
    return F ← ∪_k F_k;
```

```
Function candidate-gen(F_{k-1})
    C_k ← ∅;
    forall f₁, f₂ ∈ F_{k-1}
        with f₁ = {i₁, ... , i_{k-2}, i_{k-1}}
        and f₂ = {i₁, ... , i_{k-2}, i'_{k-1}}
        and i_{k-1} < i'_{k-1} do
        c ← {i₁, ..., i_{k-1}, i'_{k-1}};        // join f₁ and f₂
        C_k ← C_k ∪ {c};
        for each (k-1)-subset s of c do
            if (s ∉ F_{k-1}) then
                delete c from C_k;                // prune
        end
    end
    return C_k;
```

# 2   MS Apriori

sort in ascending



Figure 1: MsApriori



Figure 2: Init Pass



Figure 3: Candidate Gen MS

# 3   GSP (sequential pattern mining)



When he says is the same as the subsequence... Means that you have to remove the parenthesis and see if the elements are in the same order.

# 4 Index for Classifiers

| Metric | Description / Formula |
|---|---|
| Accuracy | $\frac{\text{Correctly Classified Examples}}{\text{Total Examples}}$ |
| Precision | $\frac{TP}{TP+FP}$: Of predicted positives, how many are correct. |
| Recall/sensitivity (TPR) | $\frac{TP}{TP+FN}$: Of actual positives, how many are identified. |
| $F_1$ Score | $2 \times \frac{p \cdot r}{p+r}$: Combines precision and recall. |
| Specificity (TNR) | $\frac{TN}{TN+FP}$: True negative rate. |
| False Positive Rate (FPR) | $1 - \text{TNR}$ |
| ROC Curve | Plot: FPR (x-axis) vs. TPR (y-axis). |

Table 1: Summary of Classification Metrics

# 5 Decision trees

**Entropy** D is the dataset. C are the different classes in the dataset.

$$\text{entropy}(D) = -\sum_{j=1}^{|C|} \Pr(C_j) \log_2 \Pr(C_j)$$

Is a positive value. The sum of the probabilities of the classes is 1. The more the data is purer, the more the entropy value is close to zero. Worst entropy is 1 (all the classes are equally distributed).

$$\text{entropy}_{A_i}(D) = \sum_{j=1}^{v} \frac{|D_j|}{|D|} \times \text{entropy}(D_j)$$

This value is the entropy if we choose to partition the data over attribute A. There are v possible values for the attribute A. Entropy($D_j$) is the entropy in the subset of data that has the value $v_j$ for the attribute A.

$$\text{gain}(D, A_i) = \text{entropy}(D) - \text{entropy}_{A_i}(D)$$

The higher the gain the better



```
.  Algorithm decisionTree(D, A, T)
1      if D contains only training examples of the same class cj ∈ C then
2          make T a leaf node labeled with class cj;
3      elseif A = ∅ then
4          make T a leaf node labeled with cj, which is the most frequent class in D
5      else   // D contains examples belonging to a mixture of classes. We select a single
6              // attribute to partition D into subsets so that each subset is purer
7          p0 = impurityEval-1(D);
8          for each attribute Ai ∈ {A1, A2, ..., Ak} do
9              pi = impurityEval-2(Ai, D)
10         end
11         Select Ag ∈ {A1, A2, ..., Ak} that gives the biggest impurity reduction,
                computed using p0 − pi;
12         if p0 − pg < threshold then    // Ag does not significantly reduce impurity p0
13             make T a leaf node labeled with cj, the most frequent class in D.
14         else                            // Ag is able to reduce impurity p0
15             Make T a decision node on Ag;
16             Let the possible values of Ag be v1, v2, ..., vm. Partition D into m
                   disjoint subsets D1, D2, ..., Dm based on the m values of Ag.
17             for each Dj in {D1, D2, ..., Dm} do
18                 if Dj ≠ ∅ then
19                     create a branch (edge) node Tj for vj as a child node of T;
20                     decisionTree(Dj, A-{Ag}, Tj)// Ag is removed
21                 end
22             end
23         end
24     end
```

# 6 Naive based classification

Product rule

$$\Pr(a1, a2) = \Pr(a1)\Pr(a2|a1)$$

Product rule for conditional probability

$$\Pr(a1, a2|c) = \Pr(a2|c)\Pr(a2|a1, c)$$

General case

$$\Pr(a_1 \cdots a_n \mid c) = \Pr(a_1 \mid c)\Pr(a_2 \mid c, a_1) \cdots \Pr(a_n \mid c, a_1 \cdots a_{n-1})$$

Conditional independence assumption

$$\Pr(a_i \mid c, a_1 \cdots a_{i-1}) = \Pr(a_i \mid c)$$

Goal:

$$\Pr(c \mid (a_1 \cdots a_n)) = \frac{\Pr(c)\Pr((a_1 \cdots a_n) \mid c)}{\Pr(a_1 \cdots a_n)}$$

$$\Pr(A_i = a_i \mid C = c_j) = \frac{n_{ij} + \lambda}{n_j + \lambda m_i}$$

Instead if there are missing values we just ignore them and use what we have

- $n_j$: # examples with $C=c_j$ in training data
- $n_{ij}$: # examples with both $A_i=a_i$ and $C=c_j$
- $m_i$: # possible values of attribute $A_i$.
- Normally, we use $\lambda = 1$

Using the law of total probability

$$\Pr(a_1 \cdots a_n) = \sum_{r=1}^{|C|} \Pr(c)\Pr((a_1 \cdots a_n) \mid c)$$

Using the product rule and the conditional independence assumption

$$\Pr(c \mid (a_1 \cdots a_n)) = \frac{\Pr(c) \prod_{i=1}^{n} \Pr(a_i \mid c)}{\sum_{r=1}^{|C|} \Pr(c) \prod_{i=1}^{n} \Pr(a_i \mid c)}$$

Adjusting the probability to account for attribute values that don't occur with that class.

# 7 Naive based text classification

$$\Pr(c|d) = \frac{\Pr(c)\Pr(d|c)}{\Pr(d)}$$

$$= \frac{\Pr(c) \prod_{k=1}^{\text{all word in d}} \Pr(w_k|c)}{\sum_{r=1}^{|C|} \Pr(c_r) \prod_{k=1}^{\text{all word in d}} \Pr(w_k|c)}$$

$$\Pr(c) = \frac{\sum_{i=1}^{|D|} \Pr(c|d_i)}{|D|}$$

$$\Pr(w_k|c) = \frac{\sum_{i=1}^{|D|} N_{ik}\Pr(c|d_i)}{\sum_{s=1}^{|V|} \sum_{i=1}^{|D|} N_{si}\Pr(c|d_i)}$$

$N_{ij}$ is the number of times that the word j appears in a document i.

# 8 K-means

**Algorithm** k-means(k, D)
1    Choose $k$ data points as the initial centroids (cluster centers)
2   **repeat**
3      **for** each data point $\mathbf{x} \in D$ do
4         compute the distance from $\mathbf{x}$ to each centroid;
5         assign $\mathbf{x}$ to the closest centroid    // a centroid represents a cluster
6      **endfor**
7      re-compute the centroids using the current cluster memberships
8   **until** the stopping criterion is met

# 9 Agglomerative clustering

**Algorithm** Agglomerative(*D*)
1   Make each data point in the data set *D* a cluster,
2   Compute all pair-wise distances of $x_1, x_2, ..., x_n \in D$;
2   **repeat**
3       find two clusters that are nearest to each other;
4       merge the two clusters form a new cluster *c*;
5       compute the distance from *c* to all other clusters;
12  **until** there is only one cluster left

Distance metrics:

- Single link: distance between two clusters is the distance between two closest data points in the two clusters.

- Complete link: distance between two clusters is the distance of two furthest data points in the two clusters.

- Average link: average distance between all points in the two clusters. (most precise metrics but also most computationally expensive)

- Centroid link: distance between the centroids of the two clusters. (doesn't consider the shape of the cluster)

# 10 LU learning

# 11 Search Engines

## 11.1 Transition Probability Matrix

Transition probability matrix:

$$A_{i,j} = \begin{cases} \frac{1}{\text{outdegree}(i)} & \text{if there is a link from } i \text{ to } j, \\ 0 & \text{otherwise.} \end{cases}$$

stochastic matrix: the sum of the elements in each row is 1.

we make stochastic by assigning an equal probability to each link for rows that are all zeros. or by removing the node without outgoing edges.

**Irreducible**   An adjacency matrix is irreducible if the graph it represents is *strongly connected*, meaning there is a path between every pair of nodes in the graph.

**Aperiodic**   A state $i$ in a Markov chain is periodic with period $k > 1$ if $k$ is the smallest number such that all paths leading from state $i$ back to state $i$ have a length that is a multiple of $k$.

A graph is *aperiodic* if there are no cycles that you cannot escape from.

**Making the Matrix Aperiodic and Irreducible**   To ensure the matrix is both aperiodic and irreducible We assign a small probability that the user will jump to a random page. (we add a damping factor)

**Final formula for page rank**

$$P = (1 - d)e + dA^T P$$

# 12 Finding holes

Use tree. Add uniformly sampled points with class Hole. Perform the regular decision tree. We add a different number of $N$ points at each node.

- The number of $N$ points for the current node $E$ is determined by the following rule (note that at the root node, the number of inherited $N$ points is 0):

1. If the number of $N$ points inherited from the parent node of $E$ is less than the number of $Y$ points in $E$, then:

    - The number of $N$ points for $E$ is increased to the number of $Y$ points in $E$.

2. Else:

    - The number of inherited $N$ points is used for $E$.

# 13 EM (expectation maximization)

Train a classifier with only the labeled documents. Use it to probabilistically classify the unlabeled documents. Use ALL the documents to train a new classifier. Iterate steps 2 and 3 to convergence.

# 14 Recommender systems

Matrix factorization Users and movies are categorized in a latent space that is composed by a fixed number of features (could be the genre, star actor). We obtain two matrixes, one for users and one for movies. The probability that an user likes the movie is the sum of product of the user and movie features.

$$P_{i,j} = U_i^T \cdot M_y$$

Where U and I are column vectors.
The learning rule is

$$u_{ki}^{t+1} = u_{ki}^t + 2\gamma(r_{ij} - p_{ij})m_{kj}^t$$

$$m_{kj}^{t+1} = m_{kj}^t + 2\gamma(r_{ij} - p_{ij})u_{ki}^t$$

# 15 Sentiment quintuple

Holder, time, entity, aspect, sentiment.

# 16 General definitions

**Continual learning**: Continual learning is the ability of a model to learn new tasks incrementally without forgetting previously learned tasks. **Class-incremental learning**: Class-incremental learning involves learning new classes incrementally while retaining the ability to classify previously learned classes. **Task-incremental learning**: Task-incremental learning involves learning new tasks sequentially, where task identity is provided during inference. **Inter-task class separation**: Inter-task class separation refers to maintaining distinct boundaries between classes from different tasks to prevent interference. **Objectives of continual learning**: The two main objectives are avoiding catastrophic forgetting and promoting knowledge transfer between tasks. **Closed world learning**: Closed world learning assumes the model only encounters data belonging to predefined classes. **Open-world learning, on-the-job learning, and continual learning after deployment**: These involve learning from data incrementally after deployment, adapting to new classes, and handling open-set scenarios. **Out-of-distribution detection**: Out-of-distribution detection identifies inputs that differ significantly from the training data's distribution. **CML (Continual Meta-Learning)**: CML enables a model to quickly adapt to new tasks using prior knowledge while mitigating forgetting. **Self-initiated open-world continual learning and adaptation**: This framework involves a model autonomously identifying and learning from novel classes in an open-world scenario