

CSE 5334 Final Preparation

* Information Retrieval Systems

- Word: a unit separated by whitespace or punctuation
- Term: a meaningful unit
- Token: a unit after tokenization
- Type: a disjoint set of tokens from documents

$$\text{Term} \div \text{Type} = \text{Word} \div \text{Token}$$

* Normalization

* Tokenization

- numbers
- stopwords
- Type construction: ① Lemmatization ② Stemming

ex. Porter, Lavis, Porter

A few rules of Porter Stemmer

Rule	Example
sses → ss	caresses → caress
ies → i	ponies → poni
ss → ss	caress → caress
s →	cats → cat

* Ranked Retrieval

① Jaccard Coefficient: from non-empty sets A, B

$$\text{Jaccard}(A, B) = \frac{|A \cap B|}{|A \cup B|}$$

- due to the Set feature, 1) X considers the freq. of terms
2) X weights the important words.
3) X implies the length of documents

② tf-idf matching score

- Term Frequency tf = # of terms in documents
- log weighted: $(1 + \log_{10}(tf)) : tf > 0$
0 : otherwise

$$\Rightarrow \sum (1 + \log_{10} tf)$$

③ TF-IDF Matrix & Cosine Similarity

- Inverse Document Frequency $idf = \log_{10} \frac{N}{df_i}$
- N : # of documents
 df_i : # of documents which contains term

- TF-IDF Matrix

term \ doc	d_1	d_2
t_1	u_{11}	v_{12}
t_2	u_{21}	v_{22}

$$\Rightarrow d_1 = [u_{11}, u_{21}]$$

$$d_2 = [v_{12}, v_{22}]$$

$$\text{Cosine Similarity} = \frac{\vec{d}_1 \cdot \vec{d}_2}{\|\vec{d}_1\| \|\vec{d}_2\|}$$

Table 5. Components of tf-idf

Term Frequency tf	Document Frequency df	Normalization
n (natural) $tf_{i,d}$	n (no)	n (none)
l (logarithmic) $1 + \log(tf_{i,d})$	1	1
a (augmented) $0.5 + 0.5 \times \frac{tf_{i,d}}{\max_{i \in d} tf_{i,d}}$	t (idf) $\log \frac{N}{df_i}$	c (cosine) $\frac{1}{\sqrt{w_1^2 + w_2^2 + \dots + w_M^2}}$
b (boolean) 1, if $tf_{i,d} > 0$, otherwise 0	p (prob idf) $\max(0, \log \frac{N - df_i}{df_i})$	u (pivoted unique) $\frac{1}{u}$
L (log average) $\frac{1 + \log tf_{i,d}}{1 + \log(\text{average } tf_{i,d})}$		b (byte size) $\frac{1}{\text{CharLength}^a} a < 1$

ex. lnc.ltn \rightarrow lnc for doc.
ltn for query

* Data Mining

- 4Vs: Volume, Velocity, Variety, Veracity

↳ quality

- Structured / Semi-Structured / Unstructured data

- Nominal — eye color, student ID

* Nom. < Ord. < Int. < Rat.

Ordinal — student grades, satisfaction level, height

Interval — Celsius & Fahrenheit Temperature, dates on calendar

Ratio — Kelvin Temperature

* Data Preprocessing

- aggregation, sampling, dimensionality reduction, feature subset selection, feature creation, discretization and binarization, attribute transformation

- random sampling
- w/o replacement sampling
- w/ replacement sampling
- Stratified sampling

* Similarity and Distance

Attribute	Similarity	Dissimilarity
Nominal	$s = \begin{cases} 1 & \text{if } p = q \\ 0 & \text{if } p \neq q \end{cases}$	$d = \begin{cases} 0 & \text{if } p = q \\ 1 & \text{if } p \neq q \end{cases}$
Ordinal	$s = 1 - \frac{ p - q }{n - 1} (s = 1 - d)$ (where n is the number of values, mapping the values to integers 0 to n-1)	$d = \frac{ p - q }{n - 1}$
Interval or Ratio	$s = -d, s = \frac{1}{1 + d}$	$d = p - q $

- Distance

$$\text{dist} = \left(\sum_{i=1}^n |p_i - q_i|^r \right)^{1/r} \begin{cases} r = 1, & \text{Manhattan distance, L1 norm} \\ r = 2, & \text{Euclidean distance, L2 norm} \\ r = \infty, & \text{Supremum (Chebyshev) distance, } L_{\max} \text{ norm, } L_{\infty} \text{ norm} \end{cases}$$

- Similarity

	SMC	Jaccard Coefficient
$M_{pq} = \begin{cases} M_{01} & p = 0, q = 1 \\ M_{10} & p = 1, q = 0 \\ M_{11} & p = 1, q = 1 \\ M_{00} & p = 0, q = 0 \end{cases}$	$\frac{M_{11} + M_{00}}{M_{11} + M_{00} + M_{10} + M_{01}}$	$\frac{M_{11}}{M_{11} + M_{10} + M_{01}}$
(The number of attributes for each case p and q) ex. p = 1000000000, q = 0000001001	$\frac{0+7}{0+7+1+2} = \frac{7}{10}$	$\frac{0}{0+1+2} = 0$

- if M_{00} has an important meaning \rightarrow SMC
- if M_{00} is meaningless \rightarrow Jaccard

$$\text{Pearson Correlation Coefficient} = \frac{\text{covariance}_{x,y}}{\sigma_x \sigma_y} = \frac{\sum_{i=1}^n (p_i - \bar{p})(q_i - \bar{q})}{\sqrt{\sum_{i=1}^n (p_i - \bar{p})^2} \sqrt{\sum_{i=1}^n (q_i - \bar{q})^2}}$$

$$\bar{p} = \frac{1}{n} \sum p_i \quad \bar{q} = \frac{1}{n} \sum q_i$$

- Overall Similarity (w/ many types of attributes)

1. Define an indicator variable σ_k for the k-th attribute as follows:

$$\sigma_k = \begin{cases} 0 & \text{(if the } k^{\text{th}} \text{ attribute is a binary asymmetric attribute AND both objects have a 0 value)} \\ & \text{or (if one of the object has a missing value for the } k^{\text{th}} \text{ attribute)} \\ 1 & \text{otherwise} \end{cases}$$

2. With the similarity of the k^{th} attribute s_k , the overall similarity between two objects is calculated as:

$$s(p, q) = \frac{\sum_{k=1}^n s_k \cdot \sigma_k}{\sum_{k=1}^n \sigma_k}$$

* Decision Tree

- Hunts

Start from a tree with a single root node containing all the training data.

Recursive:

- Check if the node is homogeneous (pure).
- If true, make the node a leaf node and label it with the class. END the branch.
- If not, continue to the next step.
- Check if the node is empty.
- If true, make the node a leaf node. END the branch.
- If not, continue to the next step.
- Check whether the node has conflicting data, a same label with different values.
- If true, mark the node as a leaf node. END the branch.
- If not, continue to the next step.
- Split the node into child nodes based on the attribute.
- During the split, the algorithm calculates the impurity of the child nodes using the 1) Gini index, 2) entropy, or 3) misclassification error.
- After splitting, the Gain is recalculated to renew the tree state.

Terminate:

- If all nodes become leaf nodes during the recursive process.
- If the split does not show certain improvement than beforehand-set threshold, regarding the impurity

- Impurity Metric

$$\text{① Misclassification Error} = 1 - \max(p_i)$$

↳ the most correct value

$$\text{ex. } 3:3:3 \rightarrow \frac{3}{9} \Rightarrow ME = \frac{6}{9}$$

$$2:3:5 \rightarrow \frac{5}{10} \Rightarrow ME = \frac{5}{10}$$

$$\text{② Gini} = 1 - \sum_{i=1}^c p_i^2 \begin{cases} c = \text{the number of classes} \\ p_i = \text{the probability of selecting an item of class } i \text{ in the node} \end{cases}$$

$$\text{cf. minimum} = 0 \quad \text{maximum} = 1 - \frac{1}{c}$$

$$\text{Gain} = \text{Gini}(\text{parent}) - \sum_{\text{child} \in \text{children}} \frac{\text{the number of data in the child}}{\text{the number of data in the parent}} \times \text{Gini}(\text{child})$$

② Entropy(t) = $-\sum_j p(j|t) \log_2 p(j|t)$ - Minimum: 0
- Maximum: $\log_2 C$

Information Gain

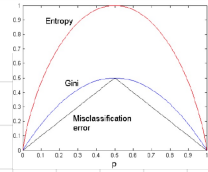
Gain_{split} = Entropy(parent) - $\sum_{\text{child} \in \text{children}} \frac{\text{the number of data in the child}}{\text{the number of data in the parent}} \times \text{Entropy}(\text{child})$

Split Info

Split Info = $-\sum_{\text{child} \in \text{children}} \frac{\text{the number of data in the child}}{\text{the number of data in the parent}} \times \log_2 \frac{\text{the number of data in the child}}{\text{the number of data in the parent}}$

Gain Ratio

Gain Ratio = $\frac{\text{Gain}_{\text{split}}}{\text{Split Info}}$



- Entropy is able to represent the maximum impurity and its derivatives along with prob. among metrics

- Pre-pruning:
 - ① Stop when the tree reaches a certain depth
 - ② Stop when the # of instances in a node is less than a certain threshold
 - ③ Stop if all instances in a node belong to the same class.
 - ④ Stop if the split does not improve impurity like Gini, Information Gain
- Post-pruning:
 - 1) Pessimistic Error
 - 2) Optimistic Error
 - 3) Reduced Error

↳ With Validation data, prune or not.

* Naive Bayes Classifier

Naive Bayes Classifier: $P(C | A_1, A_2, \dots, A_n)$

$\left\{ \begin{aligned} P(C | A_1, A_2, \dots, A_n) &= \frac{P(A_1, A_2, \dots, A_n | C) P(C)}{P(A_1, A_2, \dots, A_n)} \\ P(A_1, A_2, \dots, A_n | C) &= P(A_1 | C) P(A_2 | C) \dots P(A_n | C) \end{aligned} \right.$

$P(C_j) = \frac{\text{Number of data in class } C_j}{\text{Total number of data}}$

e.g. $P(\text{No}) = \frac{7}{10} = 0.7, P(\text{Yes}) = \frac{3}{10} = 0.3$

Discrete Attributes: $P(A_i | C_k) = \frac{\text{Number of instance having attribute } A_i}{\text{Number of Class } C_k}$

e.g. $P(\text{Status}=\text{Married} | \text{No}) = \frac{4}{7}$
 $P(\text{Refund}=\text{Yes} | \text{Yes}) = 0$

Continuous Attributes: $P(A_i | C_k) = \frac{1}{\sqrt{2\pi}\sigma^2} e^{-\frac{(x-\mu)^2}{2\sigma^2}}$

e.g. $P(\text{Income}=120 | \text{No}) = \frac{1}{\sqrt{2\pi}(54.54)^2} e^{-\frac{(120-110)^2}{2(2915)}}$

$P(X | \text{Class}=\text{No}) = P(\text{Refund}=\text{No} | \text{Class}=\text{No}) \times P(\text{Status}=\text{Married} | \text{Class}=\text{No}) \times P(\text{Income}=120K | \text{Class}=\text{No})$
 $= \frac{4}{7} \cdot \frac{4}{7} \cdot 0.0072 = 0.0024$

$P(X | \text{Class}=\text{Yes}) = P(\text{Refund}=\text{No} | \text{Class}=\text{Yes}) \times P(\text{Status}=\text{Married} | \text{Class}=\text{Yes}) \times P(\text{Income}=120K | \text{Class}=\text{Yes})$
 $= 1 \cdot 0 \cdot (1.2 \times 10^{-9}) = 0$

$P(X | \text{Class}=\text{No}) \cdot P(\text{Class}=\text{No}) = 0.0024 \cdot 0.7 = 0.0017$
 $P(X | \text{Class}=\text{Yes}) \cdot P(\text{Class}=\text{Yes}) = 0 \cdot 0.3 = 0$

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

- Prevent $P=0$ (only for Discrete Attributes)

Laplace $P(A_i | C) = \frac{\text{Number of instance having attribute } A_i + 1}{\text{Number of Class } C_k + N}$

m-estimate $P(A_i | C) = \frac{\text{Number of instance having attribute } A_i + m \cdot P(A_i)}{\text{Number of Class } C_k + m}$

N : # of classes
 m : Smoothing parameters

* KNN

- to avoid ties — 1) odd K value

2) weight factor $w = \frac{1}{\text{distance}}$

- when $K \downarrow \rightarrow$ sensitive to noise
 $K \uparrow \rightarrow$ acc. decrease.

* PEBLS

- Distance $(V_1, V_2) = \sum_{i \in C} \left| \frac{n_{1i}}{n_1} - \frac{n_{2i}}{n_2} \right| \Rightarrow$
- Weighting Method

$\delta(X, Y) = W_X W_Y \sum_{i=1}^d \text{Distance}(X_i, Y_i)^2$

Where $W_X = \frac{\text{The number of times that } X \text{ is used for prediction}}{\text{The number of times that } X \text{ is used to predict correctly}}$

* SVM

- when X_1, X_2 are given, Find $Y_2(w^T X_2 + b) \geq 1$
- non-linearly separable data: slack variable
- non-linear decision boundary: Kernel trick

* Evaluation

- Underfitting: train & test Err. \uparrow @ a model is too simple
- Overfitting: train & test \uparrow (Err.) @ a model is too complex
- Occam's Razor $\rightarrow \frac{L(M, D)}{L(M, D)} = \frac{L(M)}{L(D|M)}$
 - $L(M)$ is the length of the model
 - $L(D|M)$ is the length of the data given the model

- Missing Values

- Train
- 1) when computing Impurity
 - exclude in counting child nodes
 - include in counting parent node
 - 2) when distributing instances
 - (3 T, 1 F)

$(0T, 3F) (2T, 4F) \Rightarrow (3/4 T, 1/4 F) (6/4 T, 4/4 F)$

Test - 3) From trained tree, follow the most # of value in missing attribute

* Metrics

- Confusion Matrix

- Roc Curve

$\text{TPR} = \frac{\text{TP}}{\text{TP} + \text{FN}}$

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

draw the ROC curve by moving $\frac{1}{\text{# of positive labeled samples}}$ or $\frac{1}{\text{# of negative labeled samples}}$ threshold respectively on the each point of the ROC curve from the (1, 1) data point with the lowest threshold

- Confidence Interval

Model 1 in dataset A (size n_1 , error e_1) $\Rightarrow d_1 = e_1 - e_2$
" 2 " B (" n_2 , " e_2) $\Rightarrow d_2 = e_2 - e_3$
 $d_1 = d \pm z_{\alpha/2} \times \sigma$

* Sampling & Validation Techniques

- ① Holdout ② Random Subsampling (K Holdout) ③ Cross Valid.
- ④ Stratified Sampling ⑤ Bootstrap $k=N$: Leave-one-out

* Clustering

- Goal: Minimize Intra-Cluster distance
Maximize Inter-Cluster distance

- Partitional Clustering

$P = K! / K^k$

- Bisecting k-means: Divide the SSE cluster into 2
- SSE = $\sum_{i=1}^K \sum_{x \in C_i} \text{dist}^2(x, \mu_i)$

- Hierarchical Clustering

- Agglomerative
- Divisive — Build MST \rightarrow ① Dist. between clusters \rightarrow !

* Association Rule Mining

Rule: {Milk, Diaper} \rightarrow {Beer} $\left\{ \begin{aligned} \text{Support } s &= \frac{\sigma(\{\text{Milk, Diaper, Beer}\})}{\text{Total number of transaction}} \\ \text{Confidence } c &= \frac{\sigma(\{\text{Milk, Diaper, Beer}\})}{\sigma(\{\text{Milk, Diaper}\})} \end{aligned} \right.$

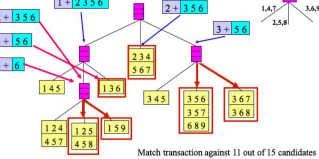
- # of Possible Association Rule: $3^n - 2^{n+1} + 1$
- # of " Itemsets : 2^n
- Apriori Algorithm
- 1) Frequent Itemset Generation (Apriori Principle)

①

1-itemsets		2-itemsets		3-itemsets	
item	support	item	support	item	support
Bread	4	Bread, Milk	3	Bread, Milk, Diaper	3
Coke	2	Bread, Beer	2	Bread, Diaper, Beer	3
Milk	4	Bread, Diaper	3	Milk, Diaper, Beer	3
Beer	3	Milk, Beer	2		
Diaper	4	Milk, Diaper	3		
Eggs	1	Beer, Diaper	3		

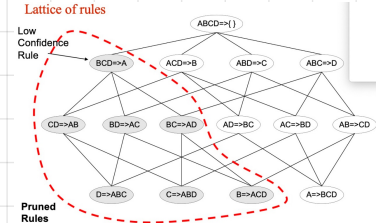
② Build Hash Tree from ①

③



2) Rule Generation

Confidence($A, B, C \rightarrow D$) \geq Confidence($A, B \rightarrow C, D$)
 \geq Confidence($A \rightarrow B, C, D$)



To solve the problem, **1)** multiple runs of the algorithm might be considered, but it is not efficient and does not guarantee the best solution with a very small probability as well. **2)** Hierarchical clustering to find initial centroids and **3)** selecting more than K initial centroids and then selecting the most widely separated among these initial centroids are another solutions. **4)** Postprocessing such as reassigning outliers or splitting and merging clusters, and **5)** bisecting K-means are also the solutions.

Handling empty clusters is another issue in the basic K-means algorithm. This could be solved by **1)** choosing the point that has the maximum distance from the centroid of the cluster with the maximum SSE and reassigning the point to the empty cluster. **2)** Repeating the algorithm is also the solution for this issue too. **3)** When it comes to running the K-means algorithm, rather than calculating the distance between all data points and all centroids after they all belong to clusters, updating the centroids after each data point is assigned to a cluster converges faster and never gets an empty cluster, but it is more expensive and introduces an order dependency. The order dependency means that the result could be different depending on which point is calculated first.

K-means algorithm itself also has limitations when the clusters of the data have different sizes, densities, or shapes. Firstly, the algorithm assumes that all clusters have the same size. However, the radius of the original clusters might be different, causing the non-optimal clusters. Secondly, the algorithm presumes that all clusters have the same density. The problem is derived from the same radius among all clusters likewise the first limitation. This leads clusters to have two dense clusters or splitted one cluster comparing from the original clusters. Lastly, the algorithm is not able to find non-globular shapes such as star-shaped clusters.