

# # 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

Term = Type ⊂ Word  
Token

## \* Normalization

## \* Tokenization

- numbers → standards
- Type construction: ① Lemmatization ② Stemming

ex. Porter, Lewis, Price

A few rules of Porter Stemmer

Rule	Example
sses → ss	caresses → caress
ies → i	ponies → pony
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 importance.  
3) X implies the length of documents

### ② tf-matching-score

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

$$\Rightarrow \sum (1 + \log tf)$$

### ③ TF-IDF Matrix & Cosine Similarity

- Inverted Document Frequency  $idf = \log \frac{N}{df_t}$  [N: # of documents  
 $df_t$ : # of documents which contains term]

$$* tf-idf = (1 + \log tf) \cdot \log \frac{N}{df_t}$$

### - TF-IDF Matrix

term	doc	df <sub>1</sub>	df <sub>2</sub>
t <sub>1</sub>	d <sub>1</sub>	u <sub>11</sub>	u <sub>21</sub>
t <sub>2</sub>	d <sub>2</sub>	u <sub>12</sub>	u <sub>22</sub>

$$\Rightarrow d_1 = [u_{11}, u_{12}] \\ d_2 = [u_{21}, u_{22}]$$

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

ex. doc. d<sub>1</sub> → doc for doc.  
d<sub>2</sub> → doc for query

Table 5. Components of tf-idf

Term Frequency tf	Document Frequency df	Normalization
n (natural)	$f_{t,d}$	
l (logarithmic)	$1 + \log(f_{t,d})$	
a (augmented)	$0.5 + 0.5 \times \frac{f_{t,d}}{\max_{t,d} f_{t,d}}$	$n(\text{no})$
b (boolean)	1, if $f_{t,d} > 0$ , otherwise 0	$t(\text{idf})$
L (log average)	$\frac{1 - \log f_{t,d}}{1 - \log(\max_{t,d} f_{t,d})}$	$\log \frac{N}{df}$
	p (prob idf)	$\max(0, \log \frac{N - df}{df})$
	u (pivot unique)	$\frac{1}{\sqrt{u_1^2 + u_2^2 + \dots + u_n^2}}$
	b (byte size)	$\frac{1}{\text{CharLength}^{1/2}} < 1$

## \* Data Mining

- AVs: Volume, Velocity, Variety, Veracity

\* (use of dimensionality

↳ the circumstance that the complexity and sparsity of the data increase as the dim. of the data increases.

- Structured / Semi-Structured / Unstructured data

- Nominal — eye color, student ID

- Ordinal — student grades, satisfaction level, height

- Interval — Celsius & Fahrenheit Temperature, dates on calendar

\* Non. C. Ord. C. Int. C. Rat.

## \* 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

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

### - Similarity

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

- if  $M_{00}$  has an important meaning → SMC
- if  $M_{00}$  is meaningless → 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  $\sigma_k$  for the k-th attribute as follows:

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

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:

1. 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.
2. Check if the node is empty.
  - If true, make the node a leaf node. END the branch.
  - If not, continue to the next step.
3. 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.
4. 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 Gini is recalculated to renew the tree state.
  - choose an attribute which maximizes a gain

Terminate:
 

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

### - Impurity Metric

- ① Misclassification Error:  $1 - \max(p_i)$

↳ the most correct value

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

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

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

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

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

$$1) \text{Entropy}(t) = -\sum_j P(j|t) \log_2 P(j|t) \quad \begin{array}{l} \text{- Minimum: } 0 \\ \text{- Maximum: } \log_2 C \end{array}$$

Information Gain

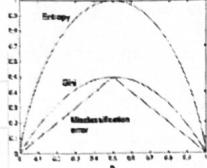
$$\text{Gain}_{\text{split}} = \text{Entropy}(\text{parent}) - \sum_{\text{child} \in \text{children}} \frac{\text{the number of data in the child}}{\text{the number of data in the parent}} \times \text{Entropy}(\text{child})$$

Split Info

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

Gain Ratio

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



- Entropy is able to represent the maximum Entropy 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:  $P(C | A_1, A_2, \dots, A_n)$

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

$X = \{\text{Refund}=\text{No}, \text{Status}=\text{Married}, \text{Income}=120K\}$

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

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

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

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

$$P(\text{Refund}=\text{Yes} | \text{Yes}) = 0$$

$P(X | \text{Class}=\text{No}) = P(\text{Refund}=\text{No} | \text{Class}=\text{No})$

$$\times P(\text{Status}=\text{Married} | \text{Class}=\text{No}) \times P(\text{Income}=120K | \text{Class}=\text{No})$$

$$= \frac{4}{7} \cdot \frac{4}{7} \cdot 0.0072 = 0.0024$$

$$P(X | \text{Class}=\text{Yes}) = P(\text{Refund}=\text{Yes} | \text{Class}=\text{Yes})$$

$$\times P(\text{Status}=\text{Married} | \text{Class}=\text{Yes}) \times P(\text{Income}=120K | \text{Class}=\text{Yes})$$

$$= 1 - 0 \cdot (1.2 \times 10^{-8}) = 0$$

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

$$\text{e.g. } P(\text{Income}=120 | \text{No}) = \frac{1}{\sqrt{2\pi(54.54)}} e^{-\frac{(120-112)^2}{2(54.54)}}$$

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

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

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

$N$ : # of classes

$m$ : Smoothing Parameters

## \* KNN

- to avoid ties — 1) odd K value

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

- when K↑ → sensitive to noise

K↑ → acc. decrease.

## \* PEBLS

$$\text{Distance}(V_1, V_2) = \sum_{j \in C} \left| \frac{n_1}{n_1} - \frac{n_2}{n_2} \right| \Rightarrow$$

Weighting Method

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

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

\* SUM + decision boundary:  $w^T b = 0$  ↗ support vector: a point which on

- when  $X_2, Y_2$  are given, Find  $Y_2(w^T X_2 + b) \geq 1$

- non-linearly separable data: slack variable

non-linear decision boundary: kernel trick

$$\begin{array}{c} \text{class} \\ \text{Yes} : 1 \quad \text{No} : 0 \end{array}$$

Distance(Single, Married) = $\frac{2}{4} - \frac{0}{4} + \frac{2}{4} - \frac{4}{4} = 1$
Distance(Single, Divorced) = $\frac{2}{4} - \frac{1}{2} + \frac{2}{4} - \frac{1}{2} = 0$
Distance(Married, Divorced) = $\frac{0}{4} - \frac{1}{2} + \frac{4}{4} - \frac{1}{2} = 1$
Distance(Refund=Yes, Refund=No) = $\frac{0}{3} - \frac{3}{7} + \frac{3}{3} - \frac{4}{7} = \frac{6}{7}$

$$\text{Margin: } \frac{2}{\sqrt{w^2}}$$

$$1 - \frac{w^T b}{\|w\|}$$

Minimize

$$\frac{1}{2} \|w\|^2$$

$$\text{Margin: } \frac{2}{\sqrt{w^2}}$$

$$1 - \frac{w^T b}{\|w\|}$$

Minimize

$$\frac{1}{2} \|w\|^2$$

## \* Evaluation

- Underfitting: train & test Err. ↑ ⚡ a model is too simple

Overfitting: train & test Err. ↑ ⚡ a model is too complex

- Occam's Razor  $L(M, D) = L(M) + L(D|M)$  ⚡ the length of the model

$L(D|M)$  is the length of the data given the model

## - Missing Values

### 1) when computing Impurity

- exclude in counting child nodes
- include in counting parent node

### 2) when distributing instances

$$(0, 1, 3, 1) \rightarrow (3/4, 1/4, 3/4)$$

$$(0, 1, 3, 1) \rightarrow (3/4, 1/4, 3/4) \rightarrow (6/16, 1/16, 6/16)$$

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

## \* Metrics

### - Confusion Matrix

$$\text{Actual} \begin{array}{c} \text{Pred.} \\ \begin{array}{ccccc} \text{TP} & \text{FP} & \text{TN} & \text{FN} \end{array} \end{array}$$

$$\text{precision} \rightarrow \frac{\text{TP}}{\text{TP} + \text{FP}}$$

### - ROC Curve

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

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

draw the ROC curve by moving  $\frac{\# \text{ of positive labeled samples}}{\# \text{ of negative labeled samples}}$  to the bottom or left

threshold respectively on the each point of the ROC curve from the (1, 1) data point with the lowest thresh

### - Confidence Interval

Model 1 in dataset A (size  $n_1$ , error  $e_1$ )  $\Rightarrow d = e_1 - e_2$

$$\therefore 2 \quad \text{in } B \quad (\sim n_2, \sim e_2) \quad d = \frac{e_1 - e_2}{\sqrt{n_2}}$$

## \* Sampling & Validation Techniques

① Holdout ② Random Subsampling (K-fold) ③ Cross Validation

④ Stratified Sampling ⑤ Bootstrap

## \* Clustering

- Goal: Minimize Intra-Cluster distance

Maximize Inter-Cluster distance

### - Partitional Clustering

$$P = K! / K^K$$

- Bisection k-means: Divide the  $\uparrow$  SSE cluster into 2

$$\text{SSE} = \sum_{i=1}^k \sum_{x \in C_i} \text{dist}^2(m_i, x)$$

MIV: can handle various shape

sensitive to noise, outliers

- MAX: less sensitive to noise, outliers able to break large clusters.

Agglomerative  $\rightarrow$  Tableau may be ok Tableau may be ok min see max.

- divisive - Build MST  $\rightarrow$   $\uparrow$  Dist. 부터 절약 clust. (Eg.)

## \* Association Rule Mining

$$\text{Support } s = \frac{\sigma(\{\text{Milk, Diaper, Beer}\})}{\text{Total number of transaction}}$$

$$\text{Rule: } \{\text{Milk, Diaper}\} \rightarrow \{\text{Beer}\}$$

$$\text{Confidence } c = \frac{\sigma(\{\text{Milk, Diaper, Beer}\})}{\sigma(\{\text{Milk, Diaper}\})}$$

- # of Possible Association rule:  $3^n - 2^{n+1} + 1$

- # of " Itemsets :  $2^n$

- Apriori Algorithm

1) Frequent Itemset Generation (Apriori Principle)

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

2) Rule Generation

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

