

## Midterm Review

### What is Data Mining

Data pre-processing: (80%: data processing; 20% modeling)

Data quality (missing values; duplicated data; etc)

#### Sampling

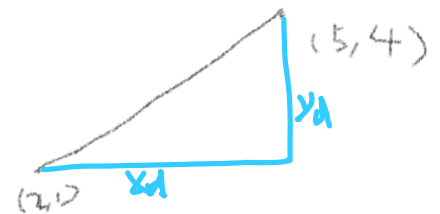
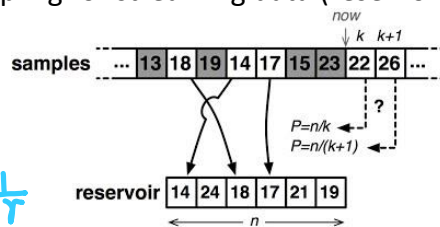
1. simple random sampling
2. sampling without replacement
3. sampling with replacement
4. stratified sampling :

(split the data into several groups, then draw random samples from each group)



ID	N	G	Y	Group	G
I	A	F	S		
II	A	F	J	E	A

#### 5. sampling for streaming data (reservoir sampling)



$$E = \sqrt{(5-2)^2 + (4-1)^2} \\ = ((5-2)^2 + (4-1)^2)^{\frac{1}{2}}$$

$$M = ((x_1 - x_2)^r + (y_1 - y_2)^r)^{\frac{1}{r}}$$

#### Similarity and dissimilarity measures

1. Euclidean distance
3. Minkowski distance ( $r=1$ ;  $r=2$ ;  $r=\infty$ ).
3. common properties of distances / similarities
4. similarity between binary vectors
- simple matching and Jaccard coefficients

$$\textcircled{1} d(x, y) \leq d(x, z) + d(z, y) \\ \textcircled{2} d(x, y) = d(y, x) \\ \textcircled{3} d(x, x) = 0$$

### Association rule mining

Concepts: set, subset, itemset, support count, support, frequent itemset

#### Association rule:

An implication expression of the form  $X \rightarrow Y$ , where  $X$  and  $Y$  are itemsets.

#### Rule evaluation metrics:

support, confidence

#### Association rule mining task:

1. Given a set of transactions, the goal of association rule mining is to find all rules have support  $\geq$  minsup threshold.
2. approaches:
  - Brute-force (very expensive)
  - Two-step approach
    - a. frequent itemset generation (very expensive)

unique items {A B C D}

frequent itemset generation strategies

a.1 Reduce the number of candidates

– Apriori principle / algorithm and its property

a.2 Reduce the number of transactions

a.3 Reduce the number of comparisons

– support counting using hash tree

b. rule generation

(generate high confidence rules from each frequent itemset)

b.1 rule generation for Apriori algorithm

b.2 factors affecting complexity of Apriori

Compact representation of frequent itemsets

Find maximal frequent itemsets

closed

---

similarity between binary vectors

simple matching and Jaccard coefficients

$$x_1 = [1, 0, 1, 1, 0, 1]$$
$$x_2 = [0, 1, 1, 0, 0, 1]$$
$$f_{11} = 1$$
$$f_{10} = 1$$
$$f_{01} = 1$$
$$f_{00} = 1$$

at least one '1' = 4

=====

Concepts: set, subset, itemset, support count, support, frequent itemset

$\{a, b, b\}$  set?  $\times$  / order  $\times$

$S_1 = \{a, b, c\}$  find subset:  $2^n = 2^3 = 8$

$\{\} = \phi$

$\{a\} \{b\} \{c\}$

$\{a, b\} \{a, c\} \{b, c\}$

$\{a, b, c\}$

K-itemset  $k=3$  is a set