# Frequent Pattern Mining

## Albert Bifet



THE UNIVERSITY OF
WAIKATO
Te Whare Wānanga o Waikato

May 2012

# COMP423A/COMP523A Data Stream Mining

**Outline**

Big Data & Real Time

# Frequent Patterns

Suppose $\mathcal{D}$ is a dataset of patterns, $t \in \mathcal{D}$, and *min_sup* is a constant.

# Frequent Patterns

Suppose $\mathcal{D}$ is a dataset of patterns, $t \in \mathcal{D}$, and *min_sup* is a constant.

### Definition
*Support* ($t$): number of patterns in $\mathcal{D}$ that are superpatterns of $t$.

# Frequent Patterns

Suppose $\mathcal{D}$ is a dataset of patterns, $t \in \mathcal{D}$, and *min_sup* is a constant.

### Definition
*Support* ($t$): number of patterns in $\mathcal{D}$ that are superpatterns of $t$.

### Definition
Pattern $t$ is *frequent* if *Support* ($t$) $\geq$ *min_sup*.

# Frequent Patterns

Suppose $\mathcal{D}$ is a dataset of patterns, $t \in \mathcal{D}$, and *min_sup* is a constant.

### Definition
*Support* ($t$): number of patterns in $\mathcal{D}$ that are superpatterns of $t$.

### Definition
Pattern $t$ is *frequent* if *Support* ($t$) $\geq$ *min_sup*.

### Frequent Subpattern Problem
Given $\mathcal{D}$ and *min_sup*, find all frequent subpatterns of patterns in $\mathcal{D}$.

# Pattern Mining

## Dataset Example

| Document | Patterns |
|:--------:|:--------:|
| d1 | abce |
| d2 | cde |
| d3 | abce |
| d4 | acde |
| d5 | abcde |
| d6 | bcd |

# Itemset Mining

| | |
|---|---|
| d1 | abce |
| d2 | cde |
| d3 | abce |
| d4 | acde |
| d5 | abcde |
| d6 | bcd |

| Support | Frequent |
|---|---|
| d1,d2,d3,d4,d5,d6 | c |
| d1,d2,d3,d4,d5 | e,ce |
| d1,d3,d4,d5 | a,ac,ae,ace |
| d1,d3,d5,d6 | b,bc |
| d2,d4,d5,d6 | d,cd |
| d1,d3,d5 | ab,abc,abe |
| | be,bce,abce |
| d2,d4,d5 | de,cde |

minimal support = 3

# Itemset Mining

| d1 | abce |
|----|------|
| d2 | cde |
| d3 | abce |
| d4 | acde |
| d5 | abcde |
| d6 | bcd |

| Support | Frequent |
|---------|----------|
| 6 | c |
| 5 | e,ce |
| 4 | a,ac,ae,ace |
| 4 | b,bc |
| 4 | d,cd |
| 3 | ab,abc,abe |
|   | be,bce,abce |
| 3 | de,cde |

# Itemset Mining

| | | Support | Frequent | Gen | Closed |
|---|---|---|---|---|---|
| d1 | abce | 6 | c | c | c |
| d2 | cde | 5 | e,ce | e | ce |
| d3 | abce | 4 | a,ac,ae,ace | a | ace |
| d4 | acde | 4 | b,bc | b | bc |
| d5 | abcde | 4 | d,cd | d | cd |
| d6 | bcd | 3 | ab,abc,abe | ab | |
| | | | be,bce,abce | be | abce |
| | | 3 | de,cde | de | cde |

# Itemset Mining

| | | Support | Frequent | Gen | Closed | Max |
|---|---|---|---|---|---|---|
| d1 | abce | 6 | c | c | c | |
| d2 | cde | 5 | e,ce | e | ce | |
| d3 | abce | 4 | a,ac,ae,ace | a | ace | |
| d4 | acde | 4 | b,bc | b | bc | |
| d5 | abcde | 4 | d,cd | d | cd | |
| d6 | bcd | 3 | ab,abc,abe | ab | | |
| | | | be,bce,abce | be | abce | abce |
| | | 3 | de,cde | de | cde | cde |

# Itemset Mining

| | |
|---|---|
| d1 | abce |
| d2 | cde |
| d3 | abce |
| d4 | acde |
| d5 | abcde |
| d6 | bcd |

| Support | Frequent | Gen | Closed | Max |
|---|---|---|---|---|
| 6 | c | c | c | |
| 5 | e,ce | e | ce | |
| 4 | a,ac,ae,ace | a | ace | |
| 4 | b,bc | b | bc | |
| 4 | d,cd | d | cd | |
| 3 | ab,abc,abe | ab | | |
| | be,bce,abce | be | abce | abce |
| 3 | de,cde | de | cde | cde |

# Itemset Mining

d1   abce
d2   cde
d3   abce
d4   acde
d5   abcde
d6   bcd

e → ce

| Support | Frequent | Gen | Closed | Max |
|---|---|---|---|---|
| 6 | c | c | c | |
| 5 | e,ce | e | ce | |
| 4 | a,ac,ae,ace | a | ace | |
| 4 | b,bc | b | bc | |
| 4 | d,cd | d | cd | |
| 3 | ab,abc,abe | ab | | |
| | be,bce,abce | be | abce | abce |
| 3 | de,cde | de | cde | cde |

# Itemset Mining

d1  ab**ce**
d2  **c**d**e**
d3  ab**ce**
d4  a**c**d**e**
d5  ab**c**d**e**
d6  bcd

| Support | Frequent | Gen | Closed | Max |
|---------|----------|-----|--------|-----|
| 6 | c | c | c | |
| 5 | e,ce | e | ce | |
| 4 | a,ac,ae,ace | a | ace | |
| 4 | b,bc | b | bc | |
| 4 | d,cd | d | cd | |
| 3 | ab,abc,abe | ab | | |
|   | be,bce,abce | be | abce | abce |
| 3 | de,cde | de | cde | cde |

# Itemset Mining

| | | Support | Frequent | Gen | Closed | Max |
|---|---|---|---|---|---|---|
| d1 | abce | 6 | c | c | c | |
| d2 | cde | 5 | e,ce | e | ce | |
| d3 | abce | 4 | a,ac,ae,ace | a | ace | |
| d4 | acde | 4 | b,bc | b | bc | |
| d5 | abcde | 4 | d,cd | d | cd | |
| d6 | bcd | 3 | ab,abc,abe | ab | | |
| | | | be,bce,abce | be | abce | abce |
| | | 3 | de,cde | de | cde | cde |

# Itemset Mining

d1   abce
d2   cde
d3   abce
d4   acde
d5   abcde
d6   bcd

$a \rightarrow ace$

| Support | Frequent | Gen | Closed | Max |
|---------|----------|-----|--------|-----|
| 6 | c | c | c | |
| 5 | e,ce | e | ce | |
| 4 | a,ac,ae,ace | a | ace | |
| 4 | b,bc | b | bc | |
| 4 | d,cd | d | cd | |
| 3 | ab,abc,abe | ab | | |
| | be,bce,abce | be | abce | abce |
| 3 | de,cde | de | cde | cde |

# Itemset Mining

| | | Support | Frequent | Gen | Closed | Max |
|---|---|---|---|---|---|---|
| d1 | abce | 6 | c | c | c | |
| d2 | cde | 5 | e,ce | e | ce | |
| d3 | abce | 4 | a,ac,ae,ace | a | ace | |
| d4 | acde | 4 | b,bc | b | bc | |
| d5 | abcde | 4 | d,cd | d | cd | |
| d6 | bcd | 3 | ab,abc,abe | ab | | |
| | | | be,bce,abce | be | abce | abce |
| | | 3 | de,cde | de | cde | cde |

# Closed Patterns

Usually, there are too many frequent patterns. We can compute a smaller set, while keeping the same information.

## Example

A set of 1000 items, has $2^{1000} \approx 10^{301}$ subsets, that is more than the number of atoms in the universe $\approx 10^{79}$

# Closed Patterns

### *A priori* property

If $t'$ is a subpattern of $t$, then $Support\,(t') \geq Support\,(t)$.

### Definition

A frequent pattern $t$ is *closed* if none of its proper superpatterns has the same support as it has.

Frequent subpatterns and their supports can be generated from closed patterns.

# Maximal Patterns

## Definition

A frequent pattern *t* is *maximal* if none of its proper superpatterns is frequent.

Frequent subpatterns can be generated from maximal patterns, but not with their support.

All maximal patterns are closed, but not all closed patterns are maximal.

# Non streaming frequent itemset miners

## Representation:

- Horizontal layout

  T1: a, b, c
  T2: b, c, e
  T3: b, d, e

- Vertical layout

  a: 1 0 0
  b: 1 1 1
  c: 1 1 0

## Search:

- Breadth-first (levelwise): Apriori
- Depth-first: Eclat, FP-Growth

# The Apriori Algorithm

APRIORI ALGORITHM

1    Initialize the item set size $k = 1$
2    Start with single element sets
3    Prune the non-frequent ones
4    **while** there are frequent item sets
5       **do** create candidates with one item more
6          Prune the non-frequent ones
7          Increment the item set size $k = k + 1$

8    Output: the frequent item sets

# The Eclat Algorithm

## Depth-First Search

- divide-and-conquer scheme : the problem is processed by splitting it into smaller subproblems, which are then processed recursively
  - **conditional database for the prefix a**
    - transactions that contain a
  - **conditional database for item sets without a**
    - transactions that not contain a
- Vertical representation
- Support counting is done by intersecting lists of transaction identifiers

# The FP-Growth Algorithm

## Depth-First Search

- ▶ divide-and-conquer scheme : the problem is processed by splitting it into smaller subproblems, which are then processed recursively
  - ▶ **conditional database for the prefix a**
    - ▶ transactions that contain a
  - ▶ **conditional database for item sets without a**
    - ▶ transactions that not contain a
- ▶ Vertical and Horizontal representation : FP-Tree
  - ▶ prefix tree with links between nodes that correspond to the same item
- ▶ Support counting is done using FP-Tree

# Mining Graph Data

### Problem
Given a data set of graphs, find frequent graphs.

| Transaction Id | Graph |
|:---:|:---:|
| 1 | O<br>$\vdots$<br>C - C - S - N<br>$\vdots$<br>O |
| 2 | O<br>$\vdots$<br>C - C - S - N<br>$\vdots$<br>C |
| 3 | N<br>$\parallel$<br>C - C - S - N |

# The gSpan Algorithm

GSPAN($g$, $D$, $min\_sup$, $S$)

Input: A graph $g$, a graph dataset $D$, $min\_sup$.
Output: The frequent graph set $S$.

1  **if** $g \neq min(g)$
2      **then return** $S$
3  insert $g$ into $S$
4  update support counter structure
5  $C \leftarrow \emptyset$
6  **for each** $g'$ that can be *right-most*
        extended from $g$ in one step
7          **do if** support($g$) $\geq$ *min_sup*
8                  **then** insert $g'$ into $C$
9  **for each** $g'$ in $C$
10          **do** $S \leftarrow$ GSPAN($g'$, $D$, $min\_sup$, $S$)
11  **return** $S$

# Mining Patterns over Data Streams

Requirements: fast, use small amount of memory and adaptive

- Type:
    - Exact
    - Approximate
- Per batch, per transaction
- Incremental, Sliding Window, Adaptive
- Frequent, Closed, Maximal patterns

# LOSSYCOUNTING

- ► Extension of LOSSYCOUNTING to Itemsets
- ► Keeps a structure with tuples $(X, \overline{freq}(X), error(X))$
- ► For each batch, to update an itemset:
  - ► Add the frequency of X in the batch to $\overline{freq}(X)$
  - ► If $\overline{freq}(X) + error(X) < bucketID$, delete this itemset
  - ► If the frequency of X in the batch in the batch is at least $\beta$, add a new tuple with $error(X) = bucketID - \beta$
- ► Uses an implementation based in :
  - ► Buffer: stores incoming transaction
  - ► Trie: forest of prefix trees
  - ► SetGen: generates itemsets supported in the current batch using apriori

# Moment

- ▶ Computes **closed** frequents itemsets in a sliding window
- ▶ Uses Closed Enumeration Tree
- ▶ Uses 4 type of Nodes:
    - ▶ Closed Nodes
    - ▶ Intermediate Nodes
    - ▶ Unpromising Gateway Nodes
    - ▶ Infrequent Gateway Nodes
- ▶ Adding transactions: closed items remains closed
- ▶ Removing transactions: infrequent items remains infrequent

# FP-Stream

- Mining Frequent Itemsets at Multiple Time Granularities
- Based in FP-Growth
- Maintains
  - pattern tree
  - tilted-time window
- Allows to answer time-sensitive queries
- Places greater information to recent data
- Drawback: time and memory complexity

# Tree and Graph Mining: Dealing with time changes

- Keep a window on recent stream elements
  - Actually, just its lattice of closed sets!
- Keep track of number of closed patterns in lattice, *N*
- Use some change detector on *N*
- When change is detected:
  - Drop stale part of the window
  - Update lattice to reflect this deletion, using deletion rule

Alternatively, sliding window of some fixed size

# Graph Coresets

## Coreset of a set *P* with respect to some problem

Small subset that approximates the original set *P*.

- ▶ Solving the problem for the coreset provides an approximate solution for the problem on *P*.

# Graph Coresets

### Coreset of a set *P* with respect to some problem

Small subset that approximates the original set *P*.

- ► Solving the problem for the coreset provides an approximate solution for the problem on *P*.

### $\delta$-tolerance Closed Graph

A graph *g* is $\delta$-*tolerance closed* if none of its proper frequent supergraphs has a weighted support $\geq (1 - \delta) \cdot$ *support(g)*.

- ► Maximal graph: 1-tolerance closed graph
- ► Closed graph: 0-tolerance closed graph.

# Graph Coresets

### Relative support of a closed graph

Support of a graph minus the relative support of its closed
supergraphs.

- ▶ The sum of the closed supergraphs' relative supports of a
  graph and its relative support is equal to its own support.

# Graph Coresets

### Relative support of a closed graph

Support of a graph minus the relative support of its closed supergraphs.

- ▶ The sum of the closed supergraphs' relative supports of a graph and its relative support is equal to its own support.

### $(s, \delta)$-coreset for the problem of computing closed graphs

Weighted multiset of frequent $\delta$-tolerance closed graphs with minimum support $s$ using their relative support as a weight.

# Graph Dataset

| Transaction Id | Graph | Weight |
|:--:|:--:|:--:|
| 1 | O<br>\|<br>C - C - S - N<br>\|<br>O | 1 |
| 2 | O<br>\|<br>C - C - S - N<br>C | 1 |
| 3 | O<br>\|<br>C - S - N<br>C | 1 |
| 4 | N<br>\|\|<br>C - C - S - N | 1 |

# Graph Coresets

| Graph | Relative Support | Support |
|---|---|---|
| C - C - S - N | 3 | 3 |
| O<br>\|<br>C - S - N | 3 | 3 |
| N<br>\|\|<br>C - S | 3 | 3 |

Table : Example of a coreset with minimum support 50% and $\delta = 1$

# Graph Coresets



Figure : Number of graphs in a $(40\%, \delta)$-coreset for NCI.

# INCGRAPHMINER

INCGRAPHMINER($D$, *min_sup*)

Input: A graph dataset $D$, and *min_sup*.
Output: The frequent graph set $G$.

1  $G \leftarrow \emptyset$
2  **for** every batch $b_t$ of graphs in $D$
3      **do** C $\leftarrow$ CORESET($b_t$, *min_sup*)
4          $G \leftarrow$ CORESET($G \cup C$, *min_sup*)
5  **return** $G$

# WINGRAPHMINER

WINGRAPHMINER($D$, $W$, $min\_sup$)

   Input: A graph dataset $D$, a size window $W$ and $min\_sup$.
   Output: The frequent graph set $G$.

```
1   G ← ∅
2   for every batch b_t of graphs in D
3       do C ← CORESET(b_t, min_sup)
4           Store C in sliding window
5           if sliding window is full
6               then R̄ ← Oldest C stored in sliding window,
                            negate all support values
7               else R̄ ← ∅
8           G ← CORESET(G ∪ C ∪ R̄, min_sup)
9   return G
```

# ADAGRAPHMINER

```
ADAGRAPHMINER(D, Mode, min_sup)
 1  G ← ∅
 2  Init ADWIN
 3  for every batch b_t of graphs in D
 4      do C ← CORESET(b_t, min_sup)
 5          R̄ ← ∅
 6          if Mode is Sliding Window
 7            then Store C in sliding window
 8                    if ADWIN detected change
 9                      then R̄ ← Batches to remove
                                   in sliding window
                                   with negative support
10          G ← CORESET(G ∪ C ∪ R̄, min_sup)
11          if Mode is Sliding Window
12            then Insert # closed graphs into ADWIN
13            else for every g in G update g's ADWIN
14  return G
```

# ADAGRAPHMINER

ADAGRAPHMINER($D$, *Mode*, *min_sup*)

```
 1   G ← ∅
 2   Init ADWIN
 3   for every batch bₜ of graphs in D
 4       do C ← CORESET(bₜ, min_sup)
 5          R̄ ← ∅
 6
 7
 8
 9


10          G ← CORESET(G ∪ C ∪ R̄, min_sup)
11
12
13          for every g in G update g's ADWIN
14   return G
```

# ADAGRAPHMINER

ADAGRAPHMINER($D$, $Mode$, $min\_sup$)

```
1   G ← ∅
2   Init ADWIN
3   for every batch b_t of graphs in D
4       do C ← CORESET(b_t, min_sup)
5           R̄ ← ∅
6           if Mode is Sliding Window
7               then Store C in sliding window
8                   if ADWIN detected change
9                       then R̄ ← Batches to remove
                                 in sliding window
                                 with negative support
10          G ← CORESET(G ∪ C ∪ R̄, min_sup)
11          if Mode is Sliding Window
12              then Insert # closed graphs into ADWIN
13
14  return G
```