WARSAW UNIVERSITY OF TECHNOLOGY
DEVELOPMENT PROGRAMME

# Finding Frequent Itemsets and Association Rules with Apriori

Marzena Kryszkiewicz

HUMAN CAPITAL
HUMAN – BEST INVESTMENT!

EUROPEAN UNION
EUROPEAN
SOCIAL FUND

Project is co-financed by European Union within European Social Fund

---

WARSAW UNIVERSITY OF TECHNOLOGY
DEVELOPMENT PROGRAMME

# Finding Frequent Itemsets

- Within each iteration $i$:
  - Determine supports of candidate itemsets of length $i$.
  - From those candidates of length $i$ that turned out frequent, create candidates of length $i + 1$.

2

---

WARSAW UNIVERSITY OF TECHNOLOGY
DEVELOPMENT PROGRAMME

# Example: Frequent 1-Itemsets

| Tid | Items |
|-----|-------|
| 1 | abce |
| 2 | abcef |
| 3 | abch |
| 4 | abe |
| 5 | acfh |
| 6 | bef |
| 7 | h |
| 8 | af |

- Let minSup = 1

- Iteration 1:
  $C_0 \rightarrow F_0$: $\varnothing_8$
  $C_1 \rightarrow F_1$: $a_6\ b_5\ c_4\ e_4\ f_4\ h_3$

3

---

WARSAW UNIVERSITY OF TECHNOLOGY
DEVELOPMENT PROGRAMME

# Example: Frequent 2-Itemsets

| Tid | Items |
|-----|-------|
| 1 | abce |
| 2 | abcef |
| 3 | abch |
| 4 | abe |
| 5 | acfh |
| 6 | bef |
| 7 | h |
| 8 | af |

- Let minSup = 1

- After iteration 1:
  $F_1$: $a_6\ b_5\ c_4\ e_4\ f_4\ h_3$

- Iteration 2:
  $C_2 \rightarrow F_2$: $ab_4\ ac_4\ ae_3\ af_3\ ah_2\ bc_3\ be_4\ bf_2\ bh_1$
  $ce_2\ cf_2\ ch_2\ ef_2\ eh_0\ fh_1$

Itemsets found as infrequent after support calculation.

4

---

WARSAW UNIVERSITY OF TECHNOLOGY
DEVELOPMENT PROGRAMME

# Example: Frequent 3-Itemsets

| Tid | Items |
|-----|-------|
| 1 | abce |
| 2 | abcef |
| 3 | abch |
| 4 | abe |
| 5 | acfh |
| 6 | bef |
| 7 | h |
| 8 | af |

- Let minSup = 1

- After iteration 2:
  $F_2$: $ab_4\ ac_4\ ae_3\ af_3\ ah_2\ bc_3\ be_4\ bf_2$
  $ce_2\ cf_2\ ch_2\ ef_2$

- Iteration 3:
  $C_3 \rightarrow F_3$: $abc_3\ abe_3\ abf_1\ abh\ ace_2\ acf_2\ ach_2$
  $aef_1\ aeh\ afh\ bce_2\ bcf_1\ bef_2\ cef_1\ ceh$
  $cfh$

Itemsets found as infrequent as supersets of infrequent itemsets.

5

---

WARSAW UNIVERSITY OF TECHNOLOGY
DEVELOPMENT PROGRAMME

# Example: Frequent 4-Itemsets

| Tid | Items |
|-----|-------|
| 1 | abce |
| 2 | abcef |
| 3 | abch |
| 4 | abe |
| 5 | acfh |
| 6 | bef |
| 7 | h |
| 8 | af |

- Let minSup = 1

- After iteration 3:
  $F_3$: $abc_3\ abe_3\ ace_2\ acf_2\ ach_2\ bce_2\ bef_2$

- Iteration 4:
  $C_4 \rightarrow F_4$: $abce_2\ acef\ aceh\ acfh$

6

## Slide 7

### Example: Frequent 5-Itemsets

| Tid | Items |
|---|---|
| 1 | abce |
| 2 | abcef |
| 3 | abch |
| 4 | abe |
| 5 | acfh |
| 6 | bef |
| 7 | h |
| 8 | af |

- Let $minSup = 1$

- After iteration 4:
  $F_4$: $abce_2$

- Iteration 5:
  $C_5$: -

7

## Slide 8

### Example: Found Frequent Itemsets

| Tid | Items |
|---|---|
| 1 | abce |
| 2 | abcef |
| 3 | abch |
| 4 | abe |
| 5 | acfh |
| 6 | bef |
| 7 | h |
| 8 | af |

$\varnothing_8$
$a_6 \ b_5 \ c_4 \ e_4 \ f_4 \ h_3$
$ab_4 \ ac_4 \ ae_3 \ af_3 \ ah_2 \ bc_3 \ be_4 \ bf_2 \ ce_2 \ cf_2 \ ch_2 \ ef_2$
$abc_3 \ abe_3 \ ace_2 \ acf_2 \ ach_2 \ bce_2 \ bef_2$
$abce_2$

- **Note:** Let $n$ be the length of a longest frequent itemset.

  Apriori finds in either $n$ or $n+1$ iterations all frequent itemsets.

8

## Slide 9

### Discovery of Association Rules with AprioriRuleGen…

- Candidate rules are built from each non-empty frequent itemset.
- Let $Z$ be a given non-empty frequent itemset. In iteration $i$, candidate rules of the form:

$$Z \setminus Y \rightarrow Y,$$

where $Y \subset Z$ and $|Y| = i$.

9

## Slide 10

### Discovery of Association Rules with AprioriRuleGen

- **Property.** Let $r_1$: $Z \setminus Y \rightarrow Y$ and $r_2$: $Z \setminus Y' \rightarrow Y'$, where $Y \subset Y'$, be association rules.
  - $conf(r_1) \geq conf(r_2)$,
  - If $conf(r_1) \leq minConf$, then $conf(r_2) \leq minConf$.

- In order to reduce the number of candidate rules, $i + 1$ item consequents of candidate rules are built from $i$ item consequents of strong association rules only.

10

## Slide 11

### Example: Discovery of **AR**s…

Frequent itemsets ($minSup = 1$):       $\varnothing_8$
$a_6 \ b_5 \ c_4 \ e_4 \ f_4 \ h_3$
$ab_4 \ ac_4 \ ae_3 \ af_3 \ ah_2 \ bc_3 \ be_4 \ bf_2 \ ce_2 \ cf_2 \ ch_2 \ ef_2$
$abc_3 \ abe_3 \ ace_2 \ acf_2 \ ach_2 \ bce_2 \ bef_2$
$abce_2$

Let $minConf = 60\%$, $Z = abce$.

**Iteration 1:**
- Consequents of candidate rules: $Y_1 = \{a, b, c, e\}$.
- Candidate rules:                              Strong association rules:
  - $bce \rightarrow a$ [2, 2/2];                $bce \rightarrow a$ [2, 2/2];
  - $ace \rightarrow b$ [2, 2/2];                $ace \rightarrow b$ [2, 2/2];
  - $abe \rightarrow c$ [2, 2/3];                $abe \rightarrow c$ [2, 2/3];
  - $abc \rightarrow e$ [2, 2/3].                $abc \rightarrow e$ [2, 2/3].

11

## Slide 12

### Example: Discovery of **AR**s…

Frequent itemsets:                      $\varnothing_8$
$a_6 \ b_5 \ c_4 \ e_4 \ f_4 \ h_3$
$ab_4 \ ac_4 \ ae_3 \ af_3 \ ah_2 \ bc_3 \ be_4 \ bf_2 \ ce_2 \ cf_2 \ ch_2 \ ef_2$
$abc_3 \ abe_3 \ ace_2 \ acf_2 \ ach_2 \ bce_2 \ bef_2$
$abce_2$

**Iteration 2 ($minConf = 60\%$, $Z = abce$):**
- Consequents of **AR**s found in iteration 1: $Y_1 = \{a, b, c, e\}$.
- Consequents of candidate rules: $Y_2 = \{ab, ac, ae, bc, be, ce\}$.
- Candidate rules:                              Strong association rules:
  - $ce \rightarrow ab$ [2, 2/2]; $ae \rightarrow bc$ [2, 2/3];        $ce \rightarrow ab$ [2, 2/2];
  - $be \rightarrow ac$ [2, 2/4]; $ac \rightarrow be$ [2, 2/4];        $bc \rightarrow ae$ [2, 2/3];
  - $bc \rightarrow ae$ [2, 2/3]; $ab \rightarrow ce$ [2, 2/4];        $ae \rightarrow bc$ [2, 2/3].

12

## Example: Discovery of **AR**s…

WARSAW UNIVERSITY OF TECHNOLOGY
DEVELOPMENT PROGRAMME

Frequent itemsets ($minSup = 1$):     $\emptyset_8$

$a_6\ b_5\ c_4\ e_4\ f_4\ h_3$

$ab_4\ ac_4\ ae_3\ af_3\ ah_2\ bc_3\ be_4\ bf_2\ ce_2\ cf_2\ ch_2\ ef_2$

$abc_3\ abe_3\ ace_2\ acf_2\ ach_2\ bce_2\ bef_2$

$abce_2$

**Iteration 3 (**$minConf = 60\%$, $Z = abce$**):**

- Consequents of **AR**s found in iteration 2: $Y_2 = \{ab, ae, bc\}$.
- Consequents of candidate rules: $Y_3 = \{abe\}$.

- Candidate rules:                    Strong association rules:
  - $c \rightarrow abe$ [2, 2/4];                              *None*

13

## Example: Found **AR**s

WARSAW UNIVERSITY OF TECHNOLOGY
DEVELOPMENT PROGRAMME

Frequent itemsets ($minSup = 1$):     $\emptyset_8$

$a_6\ b_5\ c_4\ e_4\ f_4\ h_3$

$ab_4\ ac_4\ ae_3\ af_3\ ah_2\ bc_3\ be_4\ bf_2\ ce_2\ cf_2\ ch_2\ ef_2$

$abc_3\ abe_3\ ace_2\ acf_2\ ach_2\ bce_2\ bef_2$

$abce_2$

Strong association rules ($minConf = 60\%$, $Z = abce$):

- $bce \rightarrow a$ [2, 2/2];
- $ace \rightarrow b$ [2, 2/2];
- $abe \rightarrow c$ [2, 2/3];
- $abc \rightarrow e$ [2, 2/3];
- $ce \rightarrow ab$ [2, 2/2];
- $bc \rightarrow ae$ [2, 2/3];
- $ae \rightarrow bc$ [2, 2/3].

14

## Important Operations in Apriori and AprioriRuleGen

WARSAW UNIVERSITY OF TECHNOLOGY
DEVELOPMENT PROGRAMME

- An important time-consuming operation in *Apriori* is searching *i* item candidates supported by a given transaction.
- An important time-consuming operation in *AprioriRuleGen* is searching frequent *i* itemsets (candidate rule consequents) of a given frequent itemset in order to learn their supports.
- Thus, in both cases *i* item subsets of a given itemset are searched.
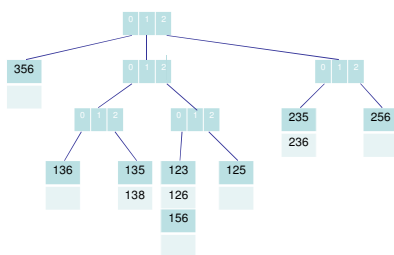
15

## Usage of a Hash Tree

WARSAW UNIVERSITY OF TECHNOLOGY
DEVELOPMENT PROGRAMME

- A hash tree is used in order to make the identification of *i* item subsets of a given itemset efficient.

- In particular, all *i* item candidate sets are stored in a hash tree.

16

## Example: Candidate 3-Itemsets in a Hash Tree

WARSAW UNIVERSITY OF TECHNOLOGY
DEVELOPMENT PROGRAMME

| itemset | represen-tation |
|---------|-----------------|
| abc | 123 |
| abe | 125 |
| abf | 126 |
| ace | 135 |
| acf | 136 |
| ach | 138 |
| aef | 156 |
| bce | 235 |
| bcf | 236 |
| bef | 256 |
| cef | 356 |



17

## Example: Searching for Subsets in a Hash Tree

WARSAW UNIVERSITY OF TECHNOLOGY
DEVELOPMENT PROGRAMME

| itemset | represen-tation |
|---------|-----------------|
| abc | 123 |
| abe | 125 |
| abf | 126 |
| … | … |

| transac-tion | represen-tation |
|--------------|-----------------|
| … | … |
| acef | 1356 |
| … | … |



- 4 subsets of transaction *acef* (after coding: 1356) has been found.

18

## Dealing with non-transactional data

- Non-transactional data are often transformed into equivalent transactional data.

19

---

## Transactional Data Set + Taxonomies

| Tid | Items |
|-----|-------|
| 1 | abce |
| 2 | abcef |
| 3 | abch |
| 4 | abe |
| 5 | acfh |
| 6 | bef |
| 7 | h |
| 8 | af |

Taxonomy tree: E (Edible) → F (Fruit), S (Sweets), V (Vitamin), M (Musical instr.); a (apple), b (banana) under F; c (chocolate) under S; e (vitamin e) under V; f (flute), h (harp) under M.

20

---

## Transactional Data Set with Taxonomies included

| Tid | Items |
|-----|-------|
| 1 | abceFSVE |
| 2 | abcefFSVME |
| 3 | abchFSME |
| 4 | abeFVE |
| 5 | acfhFSME |
| 6 | befFVME |
| 7 | hM |
| 8 | afFME |

Taxonomy tree: E (Edible) → F (Fruit), S (Sweets), V (Vitamin), M (Musical instr.); a (apple), b (banana) under F; c (chocolate) under S; e (vitamin e) under V; f (flute), h (harp) under M.

21

---

## Transactional Data Set + Taxonomies

| Tid | Items |
|-----|-------|
| 1 | abce |
| 2 | abcef |
| 3 | abch |
| 4 | abe |
| 5 | acfh |
| 6 | bef |
| 7 | h |
| 8 | af |

Taxonomy tree: E (Edible) → F (Fruit), S (Sweets), V (Vitamin), M (Musical instr.); a (apple), b (banana) under F; c (chocolate) under S; e (vitamin e) under V; f (flute), h (harp) under M.

22

---

## Transactional Data Set + Negated Items

| Tid | Items |
|-----|-------|
| 1 | abce |
| 2 | abcef |
| 3 | abch |
| 4 | abe |
| 5 | acfh |
| 6 | bef |
| 7 | h |
| 8 | af |

6 positive items {abcefh}

→

6 negated items {abcefh}

| Tid | Items |
|-----|-------|
| 1 | abcefh |
| 2 | abcefh |
| 3 | abchef |
| 4 | abecfh |
| 5 | acfhbe |
| 6 | befach |
| 7 | habcef |
| 8 | afbceh |

23

---

## Transactional Data Set + Negated Items

| Tid | Items |
|-----|-------|
| 1 | abce |
| 2 | abcef |
| 3 | abch |
| 4 | abe |
| 5 | acfh |
| 6 | bef |
| 7 | h |
| 8 | af |

6 positive items {abcefh}

→

6 negated items {abcefh}

$sup(\{befh\}) = 2$

$conf(\{bf\} \rightarrow \{eh\}) = 2/3$

$conf(\{bh\} \rightarrow \{ef\}) = 2/4$

| Tid | Items |
|-----|-------|
| 1 | abcefh |
| 2 | abcefh |
| 3 | abchef |
| 4 | abecfh |
| 5 | acfhbe |
| 6 | befach |
| 7 | habcef |
| 8 | afbceh |

24

---

## Slide 1: Transactional Data Set + Negated Items

WARSAW UNIVERSITY OF TECHNOLOGY
DEVELOPMENT PROGRAMME

### Transactional Data Set + Negated Items

| Tid | Items |
|-----|-------|
| 1 | abce*fh* |
| 2 | abce*f*h |
| 3 | abch*ef* |
| 4 | abe*cf*h |
| 5 | acfh*be* |
| 6 | bef*ach* |
| 7 | h*abcef* |
| 8 | af*bceh* |

| Tid | Items |
|-----|-------|
| 1 | 1 2 3 4 11 12 |
| 2 | 1 2 3 4 5 12 |
| 3 | 1 2 3 6 11 12 |
| 4 | 1 2 4 10 11 12 |
| 5 | 1 3 5 6 8 10 |
| 6 | 2 4 5 7 9 12 |
| 7 | 6 7 8 9 10 11 |
| 8 | 1 5 8 9 10 12 |

| Item | a | b | c | e | f | h | a | b | c | e | f | h |
|------|---|---|---|---|---|---|---|---|---|---|----|----|
| item id | 1 | 2 | 3 | 4 | 5 | 6 | 7 | 8 | 9 | 10 | 11 | 12 |

## Slide 2: Relational Data → Transactional Data

WARSAW UNIVERSITY OF TECHNOLOGY
DEVELOPMENT PROGRAMME

### Relational Data → Transactional Data

| Height | Colour | Grade |
|--------|--------|-------|
| tall | green | 5 |
| short | black | 4 |
| short | green | 4 |

| Tid | Items |
|-----|-------|
| 1 | {1, 3, 6} |
| 2 | {2, 4, 5} |
| 3 | {2, 3, 5} |

| Item | (H=tall) | (H=short) | (C=green) | (C=black) | (G=4) | (G=5) |
|------|----------|-----------|-----------|-----------|-------|-------|
| item id | 1 | 2 | 3 | 4 | 5 | 6 |
| attribute | 1 | 1 | 2 | 2 | 3 | 3 |

26

## Slide 3: Relational Data → Transactional Data

WARSAW UNIVERSITY OF TECHNOLOGY
DEVELOPMENT PROGRAMME

### Relational Data → Transactional Data

| Height | Colour | Grade |
|--------|--------|-------|
| tall | green | 5 |
| short | black | 4 |
| short | green | 4 |

| Tid | Items |
|-----|-------|
| 1 | {1, 3, 6} |
| 2 | {2, 4, 5} |
| 3 | {2, 3, 5} |

| Item | (H=tall) | (H=short) | (C=green) | (C=black) | (G=4) | (G=5) |
|------|----------|-----------|-----------|-----------|-------|-------|
| item id | 1 | 2 | 3 | 4 | 5 | 6 |
| attribute | 1 | 1 | 2 | 2 | 3 | 3 |

$\{2\} \rightarrow \{5\}$ [2, 2/2].    So, (H=short) $\rightarrow$ (G=4) [2, 100%].

$\{3\} \rightarrow \{5\}$ [1, 1/2].    So, (C=green) $\rightarrow$ (G=4) [1, 50%].

$\{2,3\} \rightarrow \{5\}$ [1, 1/1]. So, (H=short)$\wedge$(C=green) $\rightarrow$ (G=4) [1, 100%].

## Slide 4: References…

WARSAW UNIVERSITY OF TECHNOLOGY
DEVELOPMENT PROGRAMME

### References…

– Agrawal R., Imielinski T., Swami A.: Mining Associations Rules between Sets of Items in Large Databases. In: Proc. of the ACM SIGMOD Conference on Management of Data, Washington, USA (1993) 207–216

– Rakesh Agrawal, Ramakrishnan Srikant: Fast Algorithms for Mining Association Rules in Large Databases. VLDB 1994: 487-499

– Agrawal R., Mannila H., Srikant R., Toivonen H., Verkamo A.I.: Fast Discovery of Association Rules. In: Fayyad U.M., Piatetsky-Shapiro G., Smyth P., Uthurusamy R. (eds.): Advances in Knowledge Discovery and Data Mining. AAAI, CA (1996) 307–328

28

## Slide 5: References

WARSAW UNIVERSITY OF TECHNOLOGY
DEVELOPMENT PROGRAMME

### References

– Jiawei Han, Micheline Kamber, Jian Pei: Data Mining: Concept and Techniques, The Morgan Kaufmann Series in Data Management Systems, 2011

– Kryszkiewicz, M.: Concise Representations of Frequent Patterns and Association Rules. Prace Naukowe Politechniki Warszawskiej. Elektronika 142. Publishing House of the Warsaw University of Technology (2002)

– Ashok Savasere, Edward Omiecinski, Shamkant B. Navathe: An Efficient Algorithm for Mining Association Rules in Large Databases. VLDB 1995: 432-444

– Mohammed Javeed Zaki, Karam Gouda: Fast vertical mining using diffsets. KDD 2003: 326-335

29