# COMP9318: Data Warehousing and Data Mining

Assignment Project Exam Help

— L — ing —

https://eduassistpro.github.io/

Add WeChat edu_assist_pro

- **Problem definition and preliminaries**

Assignment Project Exam Help

https://eduassistpro.github.io/

Add WeChat edu_assist_pro

# What Is Association Mining?

- Association rule mining:
  - Finding frequent patterns, associations, correlations, or causal structures among sets of items or objects in transaction databases, relational databases, and other information r
  - Frequent patthttps://eduassistpro.github.io/ms, sequence, etc.) that occurs frequently in a AIS93]

    Add WeChat edu_assist_pro

- Motivation: finding regularitie
  - What products were often purchased together? — Beer and diapers?!
  - What are the subsequent purchases after buying a PC?
  - What kinds of DNA are sensitive to this new drug?
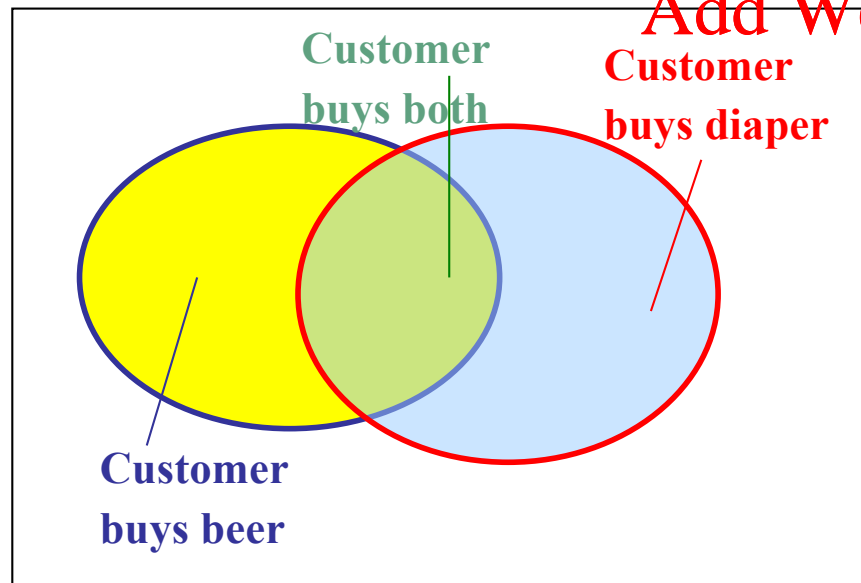  - Can we automatically classify web documents?

# Why Is Frequent Pattern or Assoiciation Mining an Essential Task in Data Mining?

- Foundation for many essential data mining tasks
  - Association, correlation, causality
  - Sequential patterns, temporal or cyclic association, partial period ~~Assignment Project Exam Help~~ media association
  - Associative c ~~https://eduassistpro.github.io/~~ alysis, iceberg cube, fascicles (semantic data co ~~Add WeChat edu_assist_pro~~ )
- Broad applications
  - Basket data analysis, cross-marketing, catalog design, sale campaign analysis
  - **Web log** (click stream) **analysis**, DNA sequence analysis, etc.

> c.f., google's spelling suggestion

# Basic Concepts: Frequent Patterns and Association Rules

| Transaction-id | Items bought |
|---|---|
| 10 | { A, B, C } |
| 20 | { A, C } |
| 30 | { A |
| 40 | { B, |

- Itemset $X=\{x_1, ..., x_k\}$
  - **Shorthand**: $x_1\ x_2\ ...\ x_k$
- Find all the rules $X \rightarrow Y$ with min confidence and support
  - support, $s$, probability that a ion contains $X \cup Y$
  - ce, $c$, conditional that a transaction lso contains $Y$.

**Customer buys both**

**Customer buys diaper**

**Customer buys beer**

*Let min_support = 50%,*
*min_conf = 70%:*
*sup(AC) = 2*
*A* ➔ *C (50%, 66.7%)*
*C* ➔ *A (50%, 100%)*

frequent itemset

association rule

# Mining Association Rules—an Example

| Transaction-id | Items bought |
|---|---|
| 10 | A, B, C |
| 20 | A, C |
| 30 | |
| 40 | |

Min. support 50%
Min. confidence 50%

Assignment Project Exam Help

https://eduassistpro.github.io/

Add WeChat edu_assist_pro

| nt pattern | Support |
|---|---|
| {A} | 75% |
| } | 50% |
| } | 50% |
| {A, C} | 50% |

For rule $A \rightarrow C$:

$\text{support} = \text{support}(\{A\} \cup \{C\}) = 50\%$

$\text{confidence} = \text{support}(\{A\} \cup \{C\})/\text{support}(\{A\}) = 66.6\%$

major computation challenge: calculate the support of itemsets
← The **frequent itemset mining** problem

- Algorithms for scalable mining of (single-dimensional Boolean) association rules in transactional databases

Assignment Project Exam Help

https://eduassistpro.github.io/

Add WeChat edu_assist_pro

# Association Rule Mining Algorithms

Candidate Generation & Verification

- Naïve algorithm
  - Enumerate all possible itemsets and check their support against $min\_sup$

  Assignment Project Exam Help

  - Generate a
    and check their confide
    against $min\_conf$

  https://eduassistpro.github.io/

  Add WeChat edu_assist_pro

- The Apriori property
  - Apriori Algorithm
  - FP-growth Algorithm

# All Candidate Itemsets for {A, B, C, D, E}

null

A    B    C    D    E

AB   AC   AD   AE   BC   BD   BE   CD   CE   DE

ABC  ABD  ABE  ACD  BCD  BCE  BDE  CDE

ABCD  ABCE  ABDE  ACDE  BCDE

ABCDE

Assignment Project Exam Help

https://eduassistpro.github.io/

Add WeChat edu_assist_pro

# Apriori Property

- A *frequent* (used to be called *large*) *itemset* is an itemset whose support is ≥ min_sup.

- Apriori property (downward closure): any subsets of a frequent ~~quent itemsets~~

- Aka the anti- of support

ABC    ABD    ~~ACD    BCD~~

ny supersets of an infrequent itemset are also infrequent itemsets"

AB    AC    AD    BC    BD    CD

A    B    C    D

# Illustrating Apriori Principle

Q: How to design an
algorithm to improve
the naïve algorithm?

null

A    B    C    D    E

Assignment Project Exam Help

AB    AC    AD    AE    BC    BD    BE    CD    CE    DE

https://eduassistpro.github.io/

Found to be
Infrequent

Add WeChat edu_assist_pro

ABC    ABD    ABE    ACD    BCD    BCE    BDE    CDE

ABCD    ABCE    ABDE    ACDE    BCDE

Pruned
supersets

ABCDE

# Apriori: A Candidate Generation-and-test Approach

- <u>Apriori pruning principle</u>: If there is any itemset which is infrequent, its superset should not be generated/tested!

Assignment Project Exam Help

- Algorithm [Ag               4]

https://eduassistpro.github.io/

  1. $C_k$ ← Perf                    idate gener (from singleton items

Add WeChat edu_assist_pro

  2. $L_k$ ← Verify $C_k$ against $L_k$
  3. $C_{k+1}$ ← generated from $L_k$
  4. Goto 2 if $C_{k+1}$ is not empty

# The Apriori Algorithm

- **Pseudo-code:**

$C_k$: Candidate itemset of size $k$
$L_k$ : frequent itemset of size $k$

```
L₁ = {frequent items};
for (k = 1; Lₖ !=∅; k++) do begin
    C_{k+1} = candidates generated from Lₖ;
    for each transaction t in database do begin
        increment the count of all candidates in C_{k+1}
        that are contained in t
    end
    L_{k+1} = candidates in C_{k+1} with min_support
end
return ∪ₖLₖ;
```

# The Apriori Algorithm—An Example

minsup = 50%

Database TDB

| Tid | Items |
|-----|-------|
| 10 | A, C, D |
| 20 | B, C, E |
| 30 | A, B, C, E |
| 40 | B, E |

$C_1$

| Itemset | sup |
|---------|-----|
| {A} | 2 |
| {B} | 3 |
| {C} | 3 |
| {D} | 1 |
| {E} | 3 |

1st scan

$L_1$

| Itemset | sup |
|---------|-----|
| {A} | 2 |
| {B} | 3 |
| {C} | 3 |
| {E} | 3 |

$C_2$

| Itemset | sup |
|---------|-----|
| {A, B} | 1 |
| {A, C} | 2 |
| {A, E} | 1 |
| {B, C} | 2 |
| {B, E} | 3 |
| {C, E} | 2 |

$C_2$

| Itemset |
|---------|
| {A, B} |
| {A, C} |
| {A, E} |
| {B, C} |
| {B, E} |
| {C, E} |

$L_2$

| Itemset | sup |
|---------|-----|
| {A, C} | 2 |
| {B, C} | 2 |
| {B, E} | 3 |
| {C, E} | 2 |

$C_3$

| Itemset |
|---------|
| {B, C, E} |

3rd scan

$L_3$

| Itemset | sup |
|---------|-----|
| {B, C, E} | 2 |

COMP9318: Data Warehousing and Data Mining

14

# Important Details of Apriori

1. How to generate candidates?

   - Step 1: self-joining $L_k$  (what's the join condition? why?)
   - Step 2: pruning

2. How to count supports of candidates?

Example of Candidate-

   - $L_3$={abc, abd, acd, ace, bcd}
   - Self-joining: $L_3*L_3$
     - abcd from abc and abd
     - acde from acd and ace
   - Pruning:
     - acde is removed because ade is not in $L_3$
   - $C_4$={abcd}

Assignment Project Exam Help

https://eduassistpro.github.io/

Add WeChat edu_assist_pro

# Generating Candidates in SQL

- Suppose the items in $L_{k-1}$ are listed in an order
- Step 1: self-joining $L_{k-1}$

insert into $C_k$

select $p.item_1,$ ... $tem_{k-1}$

from $L_{k-1}$ $p$, $L_{k-}$

where $p.item_1=q.item_1$ ... $p$ ... $em_{k-2}$, $p.item_{k-1} <$ $q.item_{k-1}$

- Step 2: pruning

forall **itemsets c in $C_k$** do

    forall **(k-1)-subsets s of c** do

        **if** *(s is not in $L_{k-1}$)* **then delete** $c$ **from** $C_k$

# Derive rules from frequent itemsets

- Frequent itemsets != association rules

- One more step is required to find association rules

- For each freq

For each prop                   t *A* of *X*,

  - Let *B* = X - *A*

  - A → B is an association rule if

    - Confidence (A → B) ≥ *min_conf*,

      where support (A → B) = support (AB), and

      confidence (A → B) = support (AB) / support (A)

# Example – deriving rules from frequent itemsets

- Suppose 234 is frequent, with supp=50%
  - Proper nonempty subsets: 23, 24, 34, 2, 3, 4, with supp=50%, 50%, 75%, 75%, 75%, 75% respectively
  - These generate these association rules:
    - 23 => 4,
    - 24 => 3,    confidence=1
    - 34 => 2,    confidence=6
    - 2 => 34,    confidence=67%
    - 3 => 24,    confidence=67%
    - 4 => 23,    confidence=67%
  - All rules have support = 50%

Assignment Project Exam Help

https://eduassistpro.github.io/

Add WeChat edu_assist_pro

$\ldots)/(N*75\%)$

Q: is there any optimization (e.g., pruning) for this step?

# Deriving rules

- To recap, in order to obtain A → B, we need to have Support(AB) and Support(A)
- This step is not as time-consuming as frequent itemsets generation
  - Why?
- It's also eas                            techniques such as parallel process
  - How?
- Do we really need candidate generation for deriving association rules?
  - Frequent-Pattern Growth (FP-Tree)

Assignment Project Exam Help

https://eduassistpro.github.io/

Add WeChat edu_assist_pro

# Bottleneck of Frequent-pattern Mining

- Multiple database scans are costly

- Mining long patterns needs many passes of scanning and andidates
  - To find fre
    - # of scans:
    - # of Candidates: $\binom{100}{1} + \binom{100}{2} + \ldots + \binom{100}{100} = 2^{100} - 1$

- Bottleneck: candidate-generation-and-test

*Can we avoid candidate generation **altogether**?*

- **FP-growth**

Assignment Project Exam Help

https://eduassistpro.github.io/

Add WeChat edu_assist_pro

# No Pain, No Gain

| | Java | Lisp | Scheme | Python | Ruby |
|---|---|---|---|---|---|
| Alice | X | | | | X |
| Bob | X | | | X | X |
| Charlie | X | | | X | X |
| Dora | | | | | |

minsup =

- Apriori:
  - L1 = {J, L, S, P, R}
  - C2 = all the $\binom{5}{2}$ combinations
    - Most of C2 do not contribute to the result
    - There is no way to tell because

# No Pain, No Gain

| | Java | Lisp | Scheme | Python | Ruby |
|---|---|---|---|---|---|
| Alice | X | | | | X |
| Bob | | | | X | X |
| Charlie | X | | | X | X |
| Dora | | | | | |

minsup =

**Ideas**:
- Keep the support set for each frequent itemset
- DFS

J ➜ JL?
J ➜ ???
Only need to look at support set for J

{A, C}

J

φ

# No Pain, No Gain

| | Java | Lisp | Scheme | Python | Ruby |
|---|---|---|---|---|---|
| Alice | X | | | | X |
| Bob | | | | X | X |
| Charlie | X | | | X | X |
| Dora | | | | | |

**Ideas**:
- Keep the support set for each frequent itemset
- DFS

minsup =

{C}

{C}

{A,C}

JPR

JP   JR

{A, C}

J

φ

...

# Notations and Invariants

- ConditonalDB:
  - DB|p = {t ∈ DB | t contains itemset p}
  - DB = DB|∅ (i.e., conditioned on nothing)
  - Shorthand:                          )
- SupportSet(p                               Set(x, DB|p)
  - {x | x mod 6 = 0 ∧ x                    =

    {x | x mod 3 = 0 ∧ x ∈ even([100]) }

- A FP-tree is equivalent to a DB|p
  - One can be converted to another
  - Next, we illustrate the alg using conditionalDB

# FP-tree Essential Idea /1

- Recursive algorithm again!

- Freq**Itemsets**(DB|p):

> easy task, as only items (not itemsets) are needed

> all frequent itemsets in DB|p belong to one of the following categories:

- X = Find**Locall** )

| patterns ~ $x_i p$ |
| patterns ~ ★$px_1$ |

output { (x p) | x ∈ X }

| patterns ~ ★$px_2$ |

*obtained via recursion*

| patterns ~ ★$px_i$ |

- Foreach x in X

  - DB*|px = GetConditionalDB$^+$(DB*|p, x)

  - 

  - Freq**Itemsets**(DB*|px)

| patterns ~ ★$px_n$ |

# No Pain, No Gain

DB|J

| | **J**ava | **L**isp | **S**cheme | **P**ython | **R**uby |
|---|---|---|---|---|---|
| Alice | X | | | | X |
| Charlie | X | | | X | X |

- Freq**Itemsets**(
  - {**P**, **R**} ← Find**Locally**Freque                    B|J)
  - Output {JP, JR}
  - Get DB*|J**P**; Freq**Itemsets**(DB*|J**P**)
  - Get DB*|J**R**; Freq**Itemsets**(DB*|J**R**)
  - // Guaranteed no other frequent itemset in DB|J

# FP-tree Essential Idea /2

- Freq**Itemsets**(DB|p):
  - If boundary con
  - X = Find**Locall**
  - [optional] DB*|p = PruneDB(

    output { (x p) | x ∈ X }
  - Foreach x in X
    - DB*|px = GetConditionalDB+(DB*|p, x)
    - [optional] if DB*|px is degenerated, then powerset(DB*|px)
    - Freq**Itemsets**(DB*|px)

Also output each item in X (appended with the conditional pattern)

Remove items not in X; potentially reduce # of transactions (∅ or dup). Improves the efficiency.

Also gets rid of items already processed before x ➜ *avoid duplicates*

# Lv 1 Recursion

- minsup = 3

| F C A M P |
|-----------|
| C B P |
| F C A M P |

DB*|P

DB*|M (sans P)

DB*|B (sans MP)

DB*|A (sans BMP)

DB*|C (sans ABMP)

DB*|F (sans CABMP)

| F C A D G I M P |
|-----------------|
| A B C F L M O |
| B F H J O W |
| B C K S P |
| A F C E L P M N |

DB

X = {F, C, A, B, M, P}

Output: F, C, A, B, M, P

| F C A M P |
|-----------|
| |
| C B P |
| F C A M P |

DB*

| F C A |
|-------|
| F C A |
| F C A |

Assignment Project Exam Help

https://eduassistpro.github.io/

Add WeChat edu_assist_pro

# Lv 2 Recursion on DB*|P

- minsup = 3

Which is actually FullDB*|CP

| F C A M P |
|-----------|
| C B P |
| F C A M P |

DB

X = {C}

Output: CP

DB*

B*|C

| C |
|---|
| C |
| C |

Context = Lv 3 recursion on DB*|CP: DB has only empty sets or X = {} ➔ immediately returns

# Lv 2 Recursion on DB*|A (sans ...)

- minsup = 3

Which is actually FullDB*|CA

| F C A |
|-------|
| F C A |
| F C A |

DB

X = {F, C}

Output: FA, CA

DB*

F C

B*|C

DB*|F

| F C |
|-----|
| F C |
| F C |

| F |
|---|
| F |
| F |

boundary case

# Different Example:
## Lv 2 Recursion on DB*|P

X = {F}

| F |
|---|
| F |

- minsup = 2

Which is actually FullDB*|AP

DB*|A

Assignment Project Exam Help

| F C |
|-----|
| F |

| F C A M P |
|-----------|
| F C B P |
| F A P |

https://eduassistpro.github.io/

DB*|C

| F |
|---|
| F |

Add WeChat edu_assist_pro

F A

B*|F

DB*

DB

X = {F, C, A}

Output: FP, CP, AP

# I will give you back the FP-tree

- An FP-tree tree of DB consists of:
  - A fixed **order** among items in DB
  - A prefix **threaded** tree of **sorted** transactions in DB
  - Header tab
- When used in the algorit put DB is always pruned (c.f., PruneDB())
  - Remove infequent items
  - Remove infequent items in every transaction

# FP-tree Example

| TID | Items bought | (ordered) frequent items |
|-----|--------------|--------------------------|
| 100 | {f, a, c, d, g, i, m, p} | {f, c, a, m, p} |
| 200 | {a, b, c, f, l, m, o} | {f, c, a, b, m} |
| 300 | {b, f, h, j, o, w} | {f, b} |
| 400 | {b, c, k, s, p} | {c, b, p} |
| 500 | {a, f, c, e, l, p, m, n} | {f, c, a, m, p} |

| TID | Items bought | (ordered) frequent items |
|---|---|---|
| **100** | *{f, a, c, d, g, i, m, p}* | *{f, c, a, m, p}* |
| **200** | *{a, b, c, f, l, m, o}* | *{f, c, a, b, m}* |
| **300** | *{b, f, h, j, o, w}* | *{f, b}* |
| **400** | *{b, c, k, s, p}* | *{c, b, p}* |
| **500** | *{a, f, c, e, l, p, m, n}* | *{f, c, a, m, p}* |

**Insert t₁**

{} → f : 1 → c : 1 → a : 1 → m : 1 → p : 1

**Insert t₂**

{} → f : 2 → c : 2 → a : 2 → m : 1, b : 1 → p : 1, m : 1

| | | |
|---|---|---|
| a | 3 | |
| | 3 | |
| m | 3 | |
| p | 3 | |

**Insert all tᵢ**

{} → f : 4, c : 1
f : 4 → c : 3, b : 1
c : 1 → b : 1 → p : 1
c : 3 → a : 3
a : 3 → m : 2, b : 1
m : 2 → p : 2
b : 1 → m : 1

| Output |
|---|
| f |
| c |
| a |
| b |
| m |
| p |

| TID | frequent items |
|-----|----------------|
| 100 | {f, c, a, m, p} |
| 200 | {f, c, a, b, m} |
| 300 | {f, b} |
| 400 | {c, b, p} |
| 500 | {f, c, a, m, p} |

**p's conditional pattern base**

| f c a   m : 2 |
|---------------|
| c   b   : 1 |

2 3 2 1 2

**Output**

pc

ed p's
tional
rn base

| C | :2 |
|---|-----|
| C | :1 |

| Item | freq | head |
|------|------|------|
| f | 4 | |
| c | 4 | |
| a | 3 | |
| b | 3 | |
| m | 3 | |
| p | 3 | |

{g

f : 4

c : 3        b : 1        b : 1

a : 3                      p : 1

m : 2    b : 1

p : 2    m : 1

Header
Table

STOP

{ }

c : 3

Assignment Project Exam Help

https://eduassistpro.github.io/

Add WeChat edu_assist_pro

| TID | frequent items |
|-----|----------------|
| 100 | {f, c, a, m, p} |
| 200 | {f, c, a, b, m} |
| 300 | {f, b} |
| 400 | {c, b, p} |
| 500 | {f, c, a, m, p} |

**m's conditional pattern base**

```
f c a   : 2
f c a b : 1
```

3 3 3 1

**Output**

| Output |
|--------|
| mf |
| mc |
| ma |

| Item | freq | head |
|------|------|------|
| f | 4 | |
| c | 4 | |
| a | 3 | |
| b | 3 | |
| m | 3 | |

{ }

f : 4          c : 1

c : 3     b : 1     b : 1

a : 3     p : 1

m : 2     b : 1

m : 1

Assignment Project Exam Help

https://eduassistpro.github.io/

Add WeChat edu_assist_pro

**Header Table**

gen_powerset

| Output |
|--------|
| mac |
| maf |
| mcf |
| macf |

{ }

f : 3

c : 3

a : 3

**b's conditional pattern base**

| f c a : 1 |
|---|
| f     : 1 |
|     c : 1 |

2 2 1

{ }

| Item | freq | head |
|---|---|---|
| f | 4 | |
| c | 4 | |
| a | 3 | |
| b | 3 | |

f : 4

c : 1

c : 3

a : 3

b : 1

STOP

Assignment Project Exam Help

https://eduassistpro.github.io/

Add WeChat edu_assist_pro

**a's conditional pattern base**

```
f c : 3
```

3 3

**Output**

| af |
| ac |

| Item | freq | head |
|------|------|------|
| f | 4 | |
| c | 4 | |
| a | 3 | |

Tree structure:

{} → f : 4 → c : 3 → a : 3

{} → c : 1

Assignment Project Exam Help

https://eduassistpro.github.io/

Add WeChat edu_assist_pro

gen_powerset

**Output**

| acf |

Header Table

{} → f : 3 → c : 3

## c's conditional pattern base

| |
|---|
| f : 3 |

| |
|---|
| 3 |

## Output

| |
|---|
| cf |

| Item | freq | head |
|------|------|------|
| f | 4 | |
| c | 4 | |

{}

f : 4

c : 3

Assignment Project Exam Help

https://eduassistpro.github.io/

Add WeChat edu_assist_pro

STOP

Header Table

{}

f : 3

STOP



Assignment Project Exam Help

https://eduassistpro.github.io/

Add WeChat edu_assist_pro

| Item | freq | head |
|------|------|------|
| f | 4 | |

f : 4

# FP-Growth vs. Apriori: Scalability With the Support Threshold

Data set T25I20D10K



Assignment Project Exam Help

https://eduassistpro.github.io/

Add WeChat edu_assist_pro

# Why Is FP-Growth the Winner?

- Divide-and-conquer:

  - decompose both the mining task and DB according to the frequent patterns obtained so far

  - leads to foc                                         databases

- Other factors

  - no candidate generation,                      te_test

  - compressed database: FP-tree structure

  - no repeated scan of entire database

  - basic ops—counting local freq items and building sub FP-tree, no pattern search and matching

Assignment Project Exam Help

https://eduassistpro.github.io/

Add WeChat edu_assist_pro