# Sequential Pattern Mining (SPM)

Sequence = Set of ordered events represented by ⟨ ⟩

Challenge → finding all subsequences

**Example:-**

| CID | TID | Transactions | | Sequential data |
|-----|-----|--------------|---|-----------------|
| 1 | 100 | a, b, c, d | CID | Sequences |
| 3 | 111 | a, f, d, e | 1 | ⟨ (abcd), (def) ⟩ |
| 1 | 122 | d, e, f | | |
| 3 | 133 | b, f, s, a | 3 | ⟨ (afde), (bfs) ⟩ |

min supp = 2

| Sid | Sequence | | |
|-----|----------|---|---|
| 10 | ⟨ a (abc) (ac) d (cf) ⟩ (1) | ⟨ (ab) c ⟩ ✓ |
| 20 | ⟨ (ad) c (bc) (ac) ⟩ | |
| 30 | ⟨ (ef) (ab) (df c b ⟩ (2) | ⟨ eg ⟩ ✗ ✓ |
| 40 | ⟨ eg (af) (cbc) ⟩ | |

## — SPM Vocabulary

- Itemset (non-empty set of items)
- Sequence (ordered list of events) ← order matters
- Event (itemset (unordered list of items))

$$(I_2 \; I_1 \; I_3)$$ where $I_1, I_2, I_3$ belongs to set of item $D$.

$e_1 = abc$          $s = ⟨ (abc)(ade) ⟩$

$e_2 \quad ade$          length = 6
                    ↓
              no. of instances of items
              in a sequence.

1 - Sequence → Length 1

Sequence with length-1 is called 1-pattern.

---

•— Subsequence $(\alpha)$ is a sequence which is part of another sequence $(\beta)$.

| Agr kuth integers exist krthy hain jisme $\alpha$ key elements $\beta$ main order main dikhte hain, tou $\alpha$ is $\beta$ ka subsequence and $\beta$ is $\alpha$ supersequence | $\alpha = \angle (ab)d >$ <br> $\beta = \angle (abc)(d.e) >$ <br> $\alpha$ subsequence of $\beta$ <br><br> $\beta = \angle (abd) > \checkmark$ |

•— Sequence database

collection of sequences and often stores as tuples

$<SID,s>$ $\rightarrow$ sequence;
↓
identifier

•— Support of a sequence in a sequence database.

• Support of sequence $(\alpha)$ in a sequence database $(S)$ is $\left(\text{no. of tuples that contains sequence}(\alpha)\right)$

calculated as no. of times sequence appears in a database.

•— Frequent Sequence.
• Sequence $(\alpha)$ is considered frequent if its support is greater or equal to specified minimum support threshold.

Algo:- Apriori based (GSP)
            Pattern-growth method (Freespan & prefixspan)
            Vertical format based mining (SPADE)
            Mining Closed sequential pattern (clospan)

Sequence pattern.
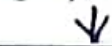   Example:-          min. supp = 2

Length of a is q &
there are 3 items in the first
3 events from I, it contributes
to Supp (a) to be $\boxed{I}$.

| SID | Sequence |
|---|---|
| 1 | < a (abc)(ac) d(ef)> |
| 2 | " |
| 3 | " |
| 4 | " |

Subsequence & support
   Example:-

< a (bc) df >     is     sub sequence    of < a (abc)(ac) d (ef)>

        ↓

It mean all key events
sequence main order key
seth aa rahay hain in

Challenges of SPMs-

• Many number of possible sequential pattern are hidden
   within large datasets. which makes (mining complex)

## Example:-    min- supp = 3

| SID | Sequence |
|-----|----------|
| 1 | < (bd) cb (ac) > |
| 2 | < (bf) (ce) b (fg) > |
| 3 | < (ah) (bdf) a bf > |
| 4 | < (be) (ce) d > |
| 5 | < a (bd) bcb (ade) > |

| Candidate | Support |
|-----------|---------|
| a | 3 |
| b | 5 |
| c | 4 |
| d | 4 |
| e | 3 |
| f | 2 X |
| g | 1 X |
| h | 1 X |

→

| candidate | Support |
|-----------|---------|
| a | 3 |
| b | 5 |
| c | 4 |
| d | 4 |
| e | 3 |

| 2-length | len (k) |
|----------|---------|
| ab | 2 |
| ac | 1 |
| ad | |
| ae | |
| (ab) | |
| | |
| (a,d) | |
| (a,e) | |

## 2) 2 events Candidate

|  | <a> | <b> | <c> | <d> | <e> |
|---|---|---|---|---|---|
| <a> | aa | ab | ac | ad | ae |
| <b> | ba | bb | bc | bd | be |
| <c> | ca | cb | c4 | cd | ce |
| <d> | da | db | dc | dd | de |
| <e> | ea | eb | ec | ed | ec |

| Candidate | Supp | Validity | Candidate | Supp | Validity |
|---|---|---|---|---|---|
| aa | 2 |  | dd | 2 |  |
| ab | 2 |  | de | 1 |  |
| ac | 1 |  | ea | 0 |  |
| ad | 2 |  | eb | 1 |  |
| ae | 1 |  | ec | 1 |  |
| ba | 3 | ✓ | ed | 1 |  |
| bb | 4 | ✓ | ee | 1 |  |
| bc | 4 | ✓ |  |  |  |
| bd | 2 |  |  |  |  |
| be | 3 | ✓ |  |  |  |
| ca | 2 |  |  |  |  |
| cb | 3 | ✓ |  |  |  |
| cc | 1 |  |  |  |  |
| cd | 2 |  |  |  |  |
| ce | 1 |  |  |  |  |
| da | 3 | ✓ |  |  |  |
| db | 3 | ✓ |  |  |  |
| dc | 2 |  |  |  |  |

# 3) Generating 1-event Candidate.

| Candidate | Supp | Validity |
|---|---|---|
| <(a,b)> | 0 | |
| <(a,c)> | 1 | |
| <(a,d)> | 1 | |
| <(b,e)> | 1 | |
| <(b,c)> | 0 | |
| <(b,d)> | 3 | ✓ |
| <(b,e)> | 1 | |
| <(c,d)> | 0 | |
| <(c,e)> | 2 | |
| <(d,e)> | 1 | |

(left margin, crossed out:)
a,b
a,c
a,d
a,e
a,f
a,g
a,h
b,c
b,d
b,e
b,f
c,b
b,h
c,d

**Frequent 2-Sequences**

2 <ba>
2 <bb>
<bc>
<be>
<cb>
<da>
<cb>
2 <(b,d)>

→ for 4th step

→ for 5th step

# 4) Generating 3-event candidate

|  | <a> | <b> | <c> | <d> | <e> |
|---|---|---|---|---|---|
| <ba> | baa | bab | bac | bad | bae |
| <bb> | bba | bbb | bbc | bbd | bbe |
| <bc> | bca | bcb | bcc | bcd | bce |
| <be> | bea | beb | bec | bed | bee |
| <cb> | cba | cbb | cbc | cbd | cbe |
| <da> | daa | dab | dac | dad | dae |
| <db> | dba | dbb | dbc | dbd | dbe |

| | | | | | | | |
|---|---|---|---|---|---|---|---|
| baa | 0 | | beb | 1 | | dba | 2 |
| bab | 1 | | bec | 0 | | dbd | 1 |
| bac | 0 | | bed | 1 | | dbe | 1 |
| bad | 0 | | bee | 0 | | | |
| bae | 0 | | cba | 2 | | | |
| bbba | 2 | | cbb | 0 | | | |
| bbb | 1 | | cbc | 1 | | | |
| bbc | 1 | | cbd | 0 | | | |
| bbd | 1 | | cbe | 1 | | | |
| bbe | 1 | | daa | 0 | | | |
| bca | 2 | | dab | 1 | | | |
| bcb | 3 ✓ | | dac | 0 | | | |
| bcc | 1 | | dad | 0 | | | |
| bcd | 2 | | dae | 0 | | | |
| bce | 1 | | dba | 2 | | | |
| bea | 0 | | dba | 1 | | | |

**5)** 2-frequent candidates.

| | a | b | c | d | e |
|---|---|---|---|---|---|
| b,d | (b,d) a | (b,d) b | (b,d) e | (b,d) d | (b,d) e |

b e b
b, d

| | | |
|---|---|---|
| (b,d) a | 3 | ✓ |
| (b,d) b | 2 | ✗ |
| (b,d) c | 2 | ✗ |
| (b,d) d | 1 | ✗ |
| (b,d) e | 1 | ✗ |