Graph and Web Mining -Motivation, Applications and Algorithms

Prof. Ehud Gudes Department of Computer Science Ben-Gurion University, Israel

Finding Sequential Patterns

Sequential Patterns Mining

Given a set of sequences, find the complete set of frequent subsequences



More Detailed Example



Motivation

- Business:
 - Customer shopping pattel
 - telephone calling patterns
 - Stock market fluctuation
 - Weblog click stream analysis
- Medical Domains:
 - Symptoms of a diseases
 - DNA sequence analysis





Definitions

- **Items**: a set of literals $\{i_1, i_2, \dots, i_m\}$
- Itemset (or event): a non-empty set of items.
- Sequence: an ordered list of itemsets, denoted as <(abc)(aef)(b)>
- A sequence <a₁...a_n> is a **subsequence** of sequence <b₁...b_m> if there exists integers i₁<...<in such that a₁ □ b_{i1},..., a_n □ b_{in}





Definitions

A <u>sequence database</u>

Seq. ID	Sequence
10	<(bd)cb(ac)>
20	<(bf)(ce)b(fg)>
30	<(ah)(bf)abf>
40	<(be)(ce)d>
50	<a(bd)bcb(ade)></a(bd)bcb(ade)>

<ad(ae)> is a <u>subsequence</u> of <<u>a(bd)bcb(ade)></u>

Given <u>support threshold</u> min_sup =2, <(bd)cb> is a <u>sequential pattern</u>



2^{*m*n*} possible candidates!



Where *m* is the number of items, and *n* in the number of transactions in the longest sequence.

 Support is the number of sequences that contain the pattern. (as in frequent itemsets, the concept of *confidence* is not defined)

 Min/Max Gap: maximum and/or minimum time gaps between adjacent elements.



Sliding Windows: consider two transactions as one as long as they are in the same time-windows.



 Multilevel: patterns that include items across different levels of hierarchy.



Multilevel





The Return of the King



The GSP Algorithm

- Developed by Srikant and Agrawal in 1996.
- Multiple-pass over the database.
- Uses generate-and-test approach.

The GSP Algorithm

- **Phase 1**: makes the first pass over database
 - To yield all the 1-element frequent sequences.
 Denoted L₁.
- **Phase 2**: the Kth pass:
 - starts with seed set found in the (k-1)th pass (L_{k-1}) to generate candidate sequences, which have one more item than a seed sequence; denoted C_k.
 - A new pass over *D* to find the support for these candidate sequences
- Phase 3: Terminates when no more frequent sequences are found

The GSP Algorithm Candidate Generation

- Joining L_{k-1} with L_{k-1}: a sequence s₁ joins with s₂ if dropping the first item from s₁ and dropping the last item from s₂ makes the same sequence.
- The added item becomes a separate event if it was a separate event in s₂, and part of the last event in s₁ otherwise.
- When joining L₁ with L₁ we need to add both ways.

Candidate Generation Example



Example

DB

SID	sequence
1	<a(abc)(ac)d(cf)></a(abc)(ac)d(cf)>
2	<(ad)c(bc)(ae)>
3	<(ef)(ab)(df)cb>
4	<eg(af)cbc></eg(af)cbc>

С	1
SEQ	Sup
<a>	4
	4
<c></c>	3
<d></d>	3
<e></e>	3
<f></f>	3
<g></g>	1



Min support =50%





Same Example – Lattice Look



GSP Drawbacks

- A huge set of candidate sequences generated.
 - Especially 2-item candidate sequence.
- Multiple Scans of database needed.
 - The length of each candidate grows by one at each database scan.
- Inefficient for mining long sequential patterns.
 - A long pattern grow up from short patterns.
 - The number of short patterns is exponential to the length of mined patterns.

The SPADE Algorithm

 SPADE (Sequential PAttern Discovery Using Equivalent Class) developed by Zaki 2001.

Ą

- A vertical format sequential pattern mining method.
- A sequence database is mapped to a large set of
 - Item: <SID, EID>
- Sequential pattern mining is performed by
 - growing the subsequences (patterns) one item at a time by Apriori candidate generation

SPADE: How It Works

Horizontal

SID	sequence
010	bequeitee
1	<a(abc)(ac)d(cf)></a(abc)(ac)d(cf)>
2	<(ad)c(bc)(ae)>
3	<(ef)(ab)(df)cb>
4	<eg(af)cbc></eg(af)cbc>

SID	EID	itemset
1	1	а
1	2	abc
1	3	ас
1	4	d
1	5	cf
2	1	ad
2	2	С
2	3	bc
2	4	ae
4	6	С

SPADE: How It Works

ID Lists for some 1-sequence

ID Lists for some 2-sequence

ā	a	ł	כ	
SI	EI	SI	EI	
			2	
	1	1	2	
1	2	2	3	
1	3	3	2	
2	1	3	5	
2	4	4	5	
3	2			
4	3			

	ab		ba				
	SID	EID(a	EID(b	SID	EID(b	EID(a	
))))	
	1	1	2	1	2	3	
•	2	1	3	2	3	4	
	3	2	5				
	4	3	5				

ID Lists for some 3-sequence

aba			•••	
SID	EID(a)	EID(b)	EID(a)	
1	1	2	3	
2	1	3	4	

SPADE: Equivalence Class



SPADE: Conclusion

- The ID Lists carry the information necessary to find support of candidates.
 Reduces scans of the sequence database.
- However, basic methodology is breadthfirst search and pruning, like GSP.

Pattern Growth: A Different Approach - PrefixSpan

- Does not require candidate generation.
- General Idea:
 - Find frequent single items.
 - Compress this information into a tree.
 - Use tree to generate a set of projected databases.
 - Each of these databases is mined separately.

Prefix and Suffix (Projection)

Let s=<a(abc)(ac)d(cf)> <a>, <aa> and <a(ab)> are prefixes of s.

Prefix	Suffix (Prefix-Based Projection)
<a>	<(abc)(ac)d(cf)>
<aa></aa>	<(_bc)(ac)d(cf)>
<ab></ab>	<(_c)(ac)d(cf)>

Mining Sequential Patterns by Prefix Projections

Step 1: find length-1 sequential patterns

<a>, , <c>, <d>, <e>, <f>

- Step 2: divide search space. The complete set of seq. pat. can be partitioned into 6 subsets:
 - The ones having prefix <a>;
 - The ones having prefix ;
 - • •
 - The ones having prefix <f>

SID	sequence
1	<a(abc)(ac)d(cf)></a(abc)(ac)d(cf)>
2	<(ad)c(bc)(ae)>
3	<(ef)(ab)(df)cb>
4	<eg(af)cbc></eg(af)cbc>

Finding Seq. Patterns with Prefix <a>

- Only need to consider projections w.r.t. <a>
 - <a>-projected database: <(abc)(ac)d(cf)>,<(_d)c(bc)(ae)>, <(_b)(df)cb>, <(_f)cbc>
- Find all the length-2 seq. pat. Having prefix <a>: <aa>, <ab>, <(ab)>, <ac>, <ad>, <af>
 - Further partition into 6 subsets
 - Having prefix <aa>;
 - • •
 - Having prefix <af>

SID	sequence
1	<a(abc)(ac)d(cf)></a(abc)(ac)d(cf)>
2	<(ad)c(bc)(ae)>
3	<(ef)(ab)(df)cb>
4	<eg(af)cbc></eg(af)cbc>

Efficiency of PrefixSpan

- No candidate sequence needs to be generated
- Projected databases keep shrinking
- Major cost of PrefixSpan: constructing projected databases
- Found to be more efficient than Spade

Constraint-Based Sequential Pattern Mining

- Constraint-based sequential pattern mining
 - Constraints: User-specified, for focused mining of desired patterns
 - How to explore efficient mining with constraints? Optimization
- Classification of constraints
 - Anti-monotone: E.g., sum(S) < 150 (If S doesn't fulfill the constraint so will super_sequence of S)
 - Monotone: E.g., count (S) > 5 (If S does fulfill the constraint so will super_sequence of S)
 - Succinct: E.g., length(S) ≥ 10, S ? (the set of sequences fullfilling the constrained can be defined precisely)
 - **Time-dependent**: E.g., min gap, max gap, total time.

Problems with Current approaches – Spade and PrefixSpan (and their variations)

- Fail (don't terminate) on database with long sequences
- Do not handle efficiently the various constraints

Our ideas - SPADE Improvement + Constraints

- Use the vertical data format
- Two phase algorithm:
 - Frequent itemset phase
 - Use the well-knows Apriori Algorithm to mine frequent itemsets.
 - Apply itemset constraints: max itemset length, items that cannot occur together.
 - Sequence phase
 - Apply sequence constraints: max gap, min gap, max/min sequence length.

Result: the CAMLS Algorithm

Frequent Itemset Phase

Use Apriori or FP-Growth to find frequent itemsets.

SID	sequence
1	<a(abc)(ac)d(cf)></a(abc)(ac)d(cf)>
2	<(ad)c(bc)(ae)>
3	<(ef)(ab)(df)cb>
4	<eg(af)cbc></eg(af)cbc>

SI D	itemse t	
1	а	
1	abc	
1	ас	
1	d	
1	cf	
2	ad	
2	С	
4	С	

 \rightarrow

Frequent Itemsets



Sequence Phase

 Similar to GSP's and SPADE's candidate generation phase – except using the frequent itemsets as seeds



So What do We Get?

- The best of both worlds:
 - Much less candidates are being generated.
 - Support check is fast.
 - Worst case: works like SPADE.
 - Tradeoff: Uses a bit more memory (for storing the frequent item-sets).

CAMLS



- Constraint-based Apriori algorithm for Mining Long Sequences
- Designed especially for efficient mining of long sequences
- Uses constraints to increase efficiency
- Outperforms both SPADE and Prefix Span on both synthetic and real data
- Appeared in DASFAA 2010

CAMLS

Makes a logical distinction between two types of constraints:

- Intra-Event: constraints that are not time related (such as items), e.g.: *Singletons*
- Inter-Event: temporal aspect of the data, i.e. values that can or cannot appear one after the other sequentially, e.g.: *Maxgap*
- Use an innovative pruning strategy

Tested Domain – predicting Machine failures

Quartz-Tungsten-Halogen lamp - used in the semiconductors industry for finding defects in a chip manufacturing process.



- Long sequences
- Restricted number of items in event
- Relatively small number of frequent events

CAMLS

- Consists of two phases corresponding to the two types of constraints:
- Event-wise: finds all frequent events satisfying intra-events constraints.
- Sequence-wise: finds all frequent sequences satisfying inter-events constraints.



Event-wise

- Iterative approach of candidate-generationand-test, based on the apriori property
- The intra-event constraints are integrated within the process, making it more efficient

CAMLS

Sequence-wise

- Iterative approach of candidate-generation-andtest, based on the apriori property
- The inter-events constraints are integrated within the process, making it more efficient
- Uses the occurrence index for efficient support calculation.
- A novel pruning strategy for redundant candidates





count and create occurrence index

Occurrence Index

- a compact representation of all occurrences of a sequence
- Structure: list of *sids*, each associated with a list of *eids*

Event-wise Example

Input minSup=2



candidates: (abc), (abd),(acd),...



No more candidates!

Sequence-wise

- $L_1 = \text{all frequent } \mathbf{1} \mathbf{sequences}$
- 2. **for** $k=2;L_{k-1}\neq \Phi;k++$ **do**
 - 1. generateCandidates(L_{k-1})
 - *2.* L_k = pruneAndSupCalc()
- 3. end for

Sequence-wise Candidate Generation

If two frequent k-sequences s' and s'' share a common k-1 prefix and s₁ is a generator, we form a new candidate

$$s' = \langle s'_{1}s'_{2...}s''_{k} \rangle \qquad s'' = \langle s''_{1}s''_{2...}s'''_{k} \rangle \\ \langle s'_{1}s'_{2...}s'_{k-1} \rangle = \langle s''_{1}s''_{2}...s''_{k-1} \rangle$$

Note that sequences grows not by one item but by all frequent events found in the first phase – i.e s''_k may be an event composed of a set of items

Sequence-wise Generator

- *maxGap* is a special kind of constraint in two ways:
 - Highly data dependant
 - Apriori property may not be applicable
- A frequent sequence that does not satisfy *maxGap* is flagged as non-generator.
- Example:
- Assume <ab> is frequent but may be pruned because the distance between a and b > maxgap
- But there are frequent sequences <ac> and <bc> and in
 <acb> all maxgap constraints are ok!
- So <ab> becomes a non-Generator but is kept in order not to prune <acb>...!

Sequence-wise Pruning: How its done

- 1. Keep a radix-ordered list of pruned sequences in current iteration
- 2. In the same iteration, one generate k-sequences from events of different size. Its possible that a k-sequence will contain another k-sequence in the same iteration.
- 3. With a new candidate:
 - 1. Check subsequence in pruned list
 - 2. Test for frequency
 - 3. Add to pruned list if needed

For example: if k-sequence <abc> was found infrequent and ksequence <a(bd)c> was generated because both b and bd are frequent, then <a(bd)c> can be pruned – this type of pruning is special to CAMLS

Original DB

Example



<aa> is added to pruned list. <a(ac)> is a super-sequence of <aa>, therefore it is pruned. <ab> does not pass maxGap,

minSup=2

maxGap=5

Event-wise

therefore it is not a generator.



Evaluation - Tested Domains

- Predicting machine failures
 - Syn-m stands for a synthetic database simulating the machine behavior with m meta-features
- Real Stocks data values
 - Rn stands for stock data (10 different stocks) for n days
- Number above rectangles indicate number of patterns found
- Note, both domains require intelligent pre-processing and discretization

CAMLS Compared with PrefixSpan



CAMLS Compared with Spade and PrefixSpan





Ben-Gurinn University of the Negry

Conclusions

 CAMLS outpeforms Spade and PrefixSpan when minSupp is low, i.e. when many sequences are generated



Thank You!