

HYPE: Hierarchical Sequential Pattern Mining

MARC PLANTEVIT, ANNE LAURENT,
MAGUELONNE TEISSEIRE

LIRMM, UNIVERSITY MONTPELLIER II, FRANCE

DOLAP'06, Arlington, November the 10th 2006

HYPE

- Introduction
- OLAP & KDD
- Sequential Patterns
- Multidimensional Framework
- Hierarchies
- Contributions
- Data Model
- Definitions
- Algorithms
- Experiments
- Conclusions and Future Work

- 1 Introduction
 - OLAP & data mining
 - Sequential Patterns
 - Multidimensional Framework
 - Hierarchies
- 2 Contributions
 - Data Model
 - Definitions
 - Algorithms
 - Experiments
- 3 Conclusions and Future Work

Introduction

OLAP & KDD
Sequential Patterns
Multidimensional Framework
Hierarchies
Contributions
Data Model
Definitions
Algorithms
Experiments
Conclusions and Future
Work

- 1 Introduction
 - OLAP & data mining
 - Sequential Patterns
 - Multidimensional Framework
 - Hierarchies

- 2 Contributions
 - Data Model
 - Definitions
 - Algorithms
 - Experiments

- 3 Conclusions and Future Work

User Navigation

- OLAP users are now decision makers.
- Users navigate in the aggregated datacube in order to discover knowledge.
 - ROLL UP, DRILL DOWN, ...

Our Goal:

Providing automatically knowledge thanks to data mining approaches

HYPE

Introduction

OLAP & KDD

Sequential Patterns

Multidimensional Framework

Hierarchies

Contributions

Data Model

Definitions

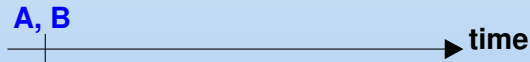
Algorithms

Experiments and Future

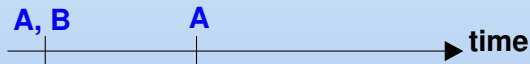
Work

- Well adapted for temporal data
- Discovering correlations between events through time.
- Several applications: marketing, decision making, protein sequence, network security, music, . . .

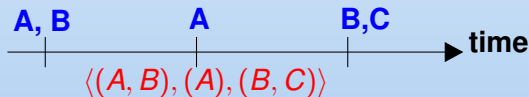
- Well adapted for temporal data
- Discovering correlations between events through time.
- Several applications: marketing, decision making, protein sequence, network security, music, . . .



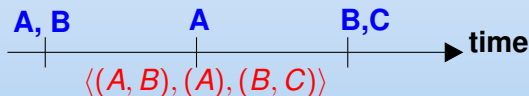
- Well adapted for temporal data
- Discovering correlations between events through time.
- Several applications: marketing, decision making, protein sequence, network security, music, . . .



- Well adapted for temporal data
- Discovering correlations between events through time.
- Several applications: marketing, decision making, protein sequence, network security, music, . . .

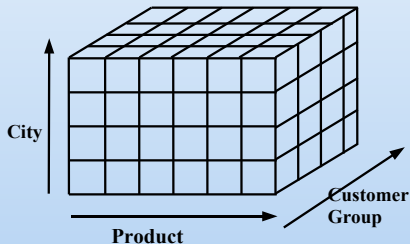


- Well adapted for temporal data
- Discovering correlations between events through time.
- Several applications: marketing, decision making, protein sequence, network security, music, ...



- ☹: Sequential patterns are quite poor (only one mined dimension)

- Knowledge are mined among one dimension: *product* dimension.
- What about the other ones ?



HYPE

Introduction

OLAP & KDD

Sequential Patterns

Multidimensional Framework

Hierarchies

Contributions

Data Model

Definitions

Algorithms

Experiments and Future

Work

- Items are not defined on one dimension, they are defined on several dimensions

HYPE

Introduction

OLAP & KDD

Sequential Patterns

Multidimensional Framework

Hierarchies

Contributions

Data Model

Definitions

Algorithms

Experiments and Future

Work

- Items are not defined on one dimension, they are defined on several dimensions
- Classical item: c

HYPE

Introduction

OLAP & KDD

Sequential Patterns

Multidimensional Framework

Hierarchies

Contributions

Data Model

Definitions

Algorithms

Experiments and Future

Work

- Items are not defined on one dimension, they are defined on several dimensions
- Classical item: c
- Multidimensional item:
(*France*, c , 100), (*Germany*, c , *)

- Items are not defined on one dimension, they are defined on several dimensions
- Classical item: c
- Multidimensional item:
 $(France, c, 100), (Germany, c, *)$
- Multidimensional sequence:

$\langle \{ (France, c, 100), (Germany, d, 54) \} \{ (*, b, 2) \} \rangle$

instead of

$\langle (c, d), b \rangle$

dilemma Support/#patterns

- Minimal support too high: too few frequent knowledge to be used and to enhance the decision making process.
- Minimal support too low: too much frequent knowledge, unusable for the decision maker.

It is very difficult to choose the right support value for mining relevant knowledge

dilemma Support/#patterns

- Minimal support too high: too few frequent knowledge to be used and to enhance the decision making process.
- Minimal support too low: too much frequent knowledge, unusable for the decision maker.

It is very difficult to choose the right support value for mining relevant knowledge

Taking hierarchies into account to solve this dilemma

- Mining rules on several levels of hierarchy.
- subsumption power.

HYPE

Introduction
 OLAP & KDD
 Sequential Patterns
 Multidimensional Framework
Hierarchies
 Contributions
 Data Model
 Definitions
 Algorithms
 Experiments
 Conclusions and Future
 Work

| | (1) | (2) | (3) |
|-----------------------|---------|--------|-----|
| Multidimensionality | No | No | Yes |
| Simulation of multi. | No | ?? | — |
| Sequential patterns | Yes | No | Yes |
| Hierarchy in patterns | Several | Single | No |

- (1) Agrawal & Srikant (1995): the pioneer approach.
- (2) Han & Fu (2001): an original approach.
- (3) Yu & Chen (2005): Using hierarchies for a smart time representation.

No approach for mining **multidimensional** sequences among **several** levels of hierarchy

HYPE

- Introduction
- OLAP & KDD
- Sequential Patterns
- Multidimensional Framework
- Hierarchies
- Contributions**
- Data Model
- Definitions
- Algorithms
- Experiments
- Conclusions and Future Work

- 1 Introduction
 - OLAP & data mining
 - Sequential Patterns
 - Multidimensional Framework
 - Hierarchies

- 2 Contributions
 - Data Model
 - Definitions
 - Algorithms
 - Experiments

- 3 Conclusions and Future Work

BLOCK :

- A database can be partitioned into different blocks according to some dimensions

| Market | Cust-Grp | Date | Place | Product |
|-----------|----------|------|---------|-----------|
| Carrefour | Educ. | 1 | Germany | beer |
| Carrefour | Educ. | 1 | Germany | pretzel |
| Carrefour | Educ. | 2 | Germany | M2 |
| Carrefour | Educ. | 3 | Germany | chocolate |
| Carrefour | Educ. | 4 | Germany | M1 |
| Carrefour | Employ. | 1 | France | soda |
| Carrefour | Employ. | 2 | France | wine |
| Carrefour | Employ. | 2 | France | chocolate |
| Carrefour | Employ. | 3 | France | M2 |
| wellmart | retir. | 1 | UK | whisky |
| wellmart | retir. | 1 | UK | pretzel |
| wellmart | retir. | 2 | UK | M2 |
| wellmart | Educ. | 1 | LA | chocolate |
| wellmart | Educ. | 2 | LA | M1 |
| wellmart | Educ. | 3 | NY | whisky |
| wellmart | Educ. | 4 | NY | soda |

$$D = \mathcal{D}_R \oplus \mathcal{D}_A \oplus \mathcal{D}_t$$

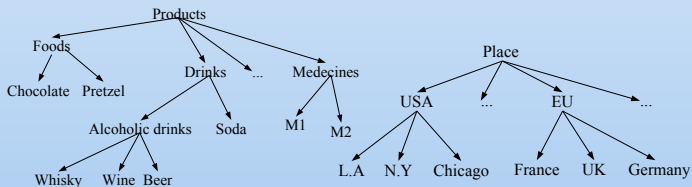
- D_t : temporal dimensions
- D_A : analysis dimensions
- D_R : reference dimensions

tuple $c = (d_1, \dots, d_n) = (r, a, t)$ where :

- r : is the restriction of c on \mathcal{D}_R
- a : is the restriction of c on \mathcal{D}_A
- t : is the restriction of c on \mathcal{D}_t

- Hierarchical relations on each analysis dimensions.
- These relations are materialized in the form of **trees** (Is-a relation).
- Only the leaves can be in the database.

Hierarchies over PLACE and PRODUCT dimensions:

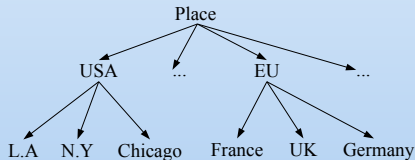


ancestor: \hat{x} is an ancestor of x according to the hierarchy.

descendant: denoted \check{x} .

$$E.U = \widehat{France}$$

$$Place = \widehat{Germany}$$



H.G. Multidimensional Item:

A tuple $e = (d_1, \dots, d_m)$ defined over the set of the analysis dimensions D_A such that $d_i \in \{label(T_i)\}$.

Examples : $(France, Chocolate)$,
 $(Germany, Drinks)$

H.G. Multidimensional Item:

A tuple $e = (d_1, \dots, d_m)$ defined over the set of the analysis dimensions D_A such that $d_i \in \{label(T_i)\}$.

Examples : $(France, Chocolate)$,
 $(Germany, Drinks)$

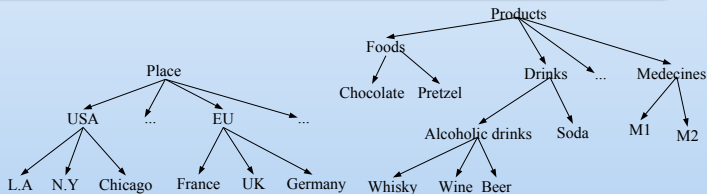
Hierarchical Inclusion

Let $e = (d_1, \dots, d_m)$ and $e' = (d'_1, \dots, d'_m)$, then:

- e is **more general** than e' ($e >_h e'$) if
 $\forall d_i, d_i = \hat{d}'_i$ or $d_i = d'_i$
- e is **more specific** than e' ($e <_h e'$) if
 $\forall d_i, d_i = \check{d}'_i$ or $d_i = d'_i$
- e and e' are **incomparable** if there is no relation between them ($e \not>_h e'$ et $e' \not>_h e$)

Example: Relation between items

- $(USA, Drinks) >_h (USA, soda)$.
- $(France, wine) <_h (EU, Alcoholic drinks)$.
- $(France, wine)$ and $(USA, soda)$ are incomparable.



H.G. Multidimensional Itemset:

$i = \{e_1, \dots, e_k\}$ where all items are all incomparable.

$\{(France, wine), (USA, soda)\}$ YES

$\{(France, wine), (EU, Alcohol.D.)\}$ NO:
 $(France, wine) <_h (EU, Alcohol.D.)$

H.G. Multidimensional Sequence

$s = \langle i_1, \dots, i_j \rangle$ is an ordered list of h-generalized itemsets.

$\langle \{(France, wine), (USA, soda)\} \{ (Germany, beer) \} \rangle$

Item Supported by a tuple

A tuple b supports an item e if $\Pi_{D_A}(b) <_h e$.

Tuple (*Carrefour, Educ, 1, France, wine*) supports the item (*EU, Alcohol.D.*).

Sequence Supported by a Block

A block supports a sequence $\langle i_1, \dots, i_l \rangle$ if
 $\forall j = 1 \dots l, \exists d_j \in \text{Dom}(D_j)$, for each item e from i_j ,
 $\exists t = (r, e, d_j)$ or $t = (r, \check{e}, d_j) \in T$ w.r.t
 $d_1 < d_2 < \dots < d_l$.

Let :

- D_R : the reference dimensions
- DB : the database partitioned into a set of block B_{T,D_R}
- A sequence S

Support of S

$$\text{support}(S) = \frac{|\{B \in B_{DB,D_R} \text{ s.t. } B \text{ supports } S\}|}{|B_{DB,D_R}|}$$

$D_R = \{Market, Cust - Grp\}$, $D_A = \{Place, Product\}$
et $D_T = \{Date\}$, $minsupp = 2$

Let us search the support of the sequence
 $S = \langle \{(UE, pretzel), (UE, Alcoholic Drinks)\}$
 $\{(UE, M2)\}\rangle$

| Market | Cust-grp | Date | Place | Product |
|-----------|----------|------|---------|-----------|
| Carrefour | Educ. | 1 | Germany | beer |
| Carrefour | Educ. | 1 | Germany | pretzel |
| Carrefour | Educ. | 2 | Germany | M2 |
| Carrefour | Educ. | 3 | Germany | chocolate |
| Carrefour | Educ. | 4 | Germany | M1 |
| Carrefour | Employ. | 1 | France | soda |
| Carrefour | Employ. | 2 | France | wine |
| Carrefour | Employ. | 2 | France | chocolate |
| Carrefour | Employ. | 3 | France | M2 |
| wellmart | retir. | 1 | UK | whisky |
| wellmart | retir. | 1 | UK | pretzel |
| wellmart | retir. | 2 | UK | M2 |
| wellmart | Educ. | 1 | LA | chocolate |
| wellmart | Educ. | 2 | LA | M1 |
| wellmart | Educ. | 3 | NY | whisky |
| wellmart | Educ. | 4 | NY | soda |

HYPE

- Introduction
- OLAP & KDD
- Sequential Patterns
- Multidimensional Framework
- Hierarchies
- Contributions
- Data Model
- Definitions
- Algorithms
- Experiments
- Conclusions and Future Work

Block 1

| | | | | |
|-----------|-------|---|----------------|----------------|
| Carrefour | Educ. | 1 | Germany | pretzel |
| Carrefour | Educ. | 1 | Germany | beer |
| Carrefour | Educ. | 2 | Germany | M2 |
| Carrefour | Educ. | 3 | Germany | chocolate |
| Carrefour | Educ. | 4 | Germany | M1 |

Block 1 supports S : $support(S) ++$

Block 2

| | | | | |
|-----------|---------|---|---------------|----------------|
| Carrefour | Employ. | 1 | France | soda |
| Carrefour | Employ. | 2 | France | pretzel |
| Carrefour | Employ. | 2 | France | wine |
| Carrefour | Employ. | 3 | France | M2 |

Block 2 supports S : $support(S) ++$

Block 3

| | | | | |
|----------|--------|---|-----------|----------------|
| wellmart | retir. | 1 | UK | pretzel |
| wellmart | retir. | 1 | UK | whisky |
| wellmart | retir. | 2 | UK | M2 |

Block 3 supports S : $support(S) + +$

Block 4

| | | | | |
|----------|-------|---|----|-----------|
| wellmart | Educ. | 1 | LA | chocolate |
| wellmart | Educ. | 2 | LA | M1 |
| wellmart | Educ. | 3 | NY | whisky |
| wellmart | Educ. | 4 | NY | soda |

Block 4 does not support S

Block 3

| | | | | |
|----------|--------|---|----|---------|
| wellmart | retir. | 1 | UK | pretzel |
| wellmart | retir. | 1 | UK | whisky |
| wellmart | retir. | 2 | UK | M2 |

Block 3 supports S : $support(S) + +$

Block 4

| | | | | |
|----------|-------|---|----|-----------|
| wellmart | Educ. | 1 | LA | chocolate |
| wellmart | Educ. | 2 | LA | M1 |
| wellmart | Educ. | 3 | NY | whisky |
| wellmart | Educ. | 4 | NY | soda |

Block 4 does not support S

$support(S) = 3 \geq minsupp$

- s is frequent

HYPE

Introduction

OLAP & KDD

Sequential Patterns

Multidimensional Framework

Hierarchies

Contributions

Data Model

Definitions

Algorithms

Experiments and Future

Work

Frequent Item Generation

- to mine all the **maximally specified items**.
- levelwise generation

Frequent Sequence Generation

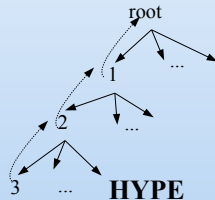
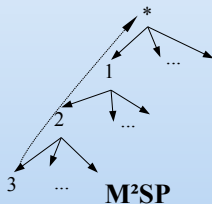
- anti-monotonicity property of the support
- Apriori-like method (generate and prune)
- use of prefix tree to store the sequences (PSP)

- required: data pre-processing
(group by *date*, D_1, \dots, D_n)
- Support Sequence Count:
SupportCount(s, DB, \mathcal{D}_R)
 - For each block : **SupportBlock**(s, B)
- **anchoring** operation ($\sigma_{condition}(B) \mapsto C'$ with $B' \subseteq B$)

complexity

- n_B : # tuples in B
- $m = |D_A|$: # analysis dimensions
- P_{max} : maximal depth of the hierarchies
- SupportBlock: $O(P_{max} \times n_B \times m \times \log n_B)$
- SupportCount: $O(I \times P_{max} \times n_{max} \times m \times \log n_{max})$

- (M²SP): a "binary" management of wild-card value.
- *HYPE*: A better management thanks to the taking hierarchies into account.
- More Accurate knowledge.



HYPE

Introduction

OLAP & KDD

Sequential Patterns

Multidimensional Framework

Hierarchies

Contributions

Data Model

Definitions

Algorithms

Experiments

Conclusions and Future

Work

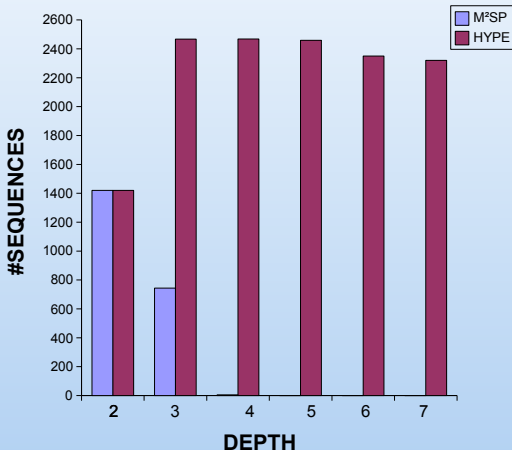
- Synthetic database
- 10,000 tuples
- $|D_A| = 5$
- Studying # frequent sequences according to:
 - Hierarchy depth (specialization degree)
 - user-defined minimal threshold
- Comparison with $M^2SP(-\alpha)$:

frequent sequences over Hierarchy depth:

- minsup=0.3, average degree = 3

HYPE

Introduction
 OLAP & KDD
 Sequential Patterns
 Multidimensional Framework
 Hierarchies
 Contributions
 Data Model
 Definitions
 Algorithms
 Experiments
 Conclusions and Future
 Work

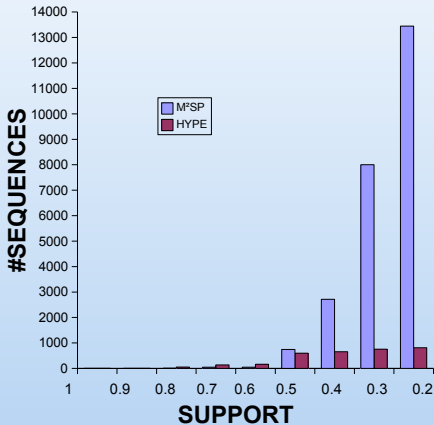


frequent sequences over minimal support:

- Dense Database (lower degree)

HYPE

- Introduction
- OLAP & KDD
- Sequential Patterns
- Multidimensional Framework
- Hierarchies
- Contributions
- Data Model
- Definitions
- Algorithms
- Experiments
- Conclusions and Future Work

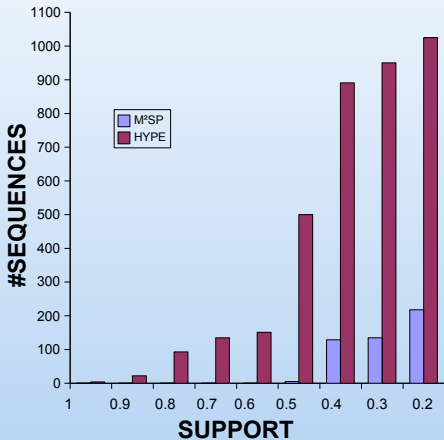


frequent sequences over minimal support:

- Sparse Database (higher degree)

HYPE

[Introduction](#)
[OLAP & KDD](#)
[Sequential Patterns](#)
[Multidimensional Framework](#)
[Hierarchies](#)
Contributions
[Data Model](#)
[Definitions](#)
[Algorithms](#)
Experiments
[Conclusions](#) and [Future Work](#)



HYPE

Introduction

OLAP & KDD

Sequential Patterns

Multidimensional Framework

Hierarchies

Contributions

Data Model

Definitions

Algorithms

Experiments

Conclusions and Future Work

- 1 Introduction
 - OLAP & data mining
 - Sequential Patterns
 - Multidimensional Framework
 - Hierarchies

- 2 Contributions
 - Data Model
 - Definitions
 - Algorithms
 - Experiments

- 3 Conclusions and Future Work

Managing hierarchies in multidimensional sequential pattern mining

- H.G. multidimensional sequential patterns.
- More accurate knowledge.
- An approach more efficient according to "min. support/mined knowledge dilemma" thanks to subsumption ability.

| | Agrawal | Jiawei Han | Yu | HYPE |
|-----------------------|---------|------------|-----|------|
| Multidimensionality | No | No | Yes | Yes |
| Simulation of multi. | No | ?? | — | Yes |
| Sequential patterns | Yes | No | Yes | Yes |
| Hierarchy in patterns | Several | Single | No | Yes |

Future Work

- Use of *condensed* representation (closed, free patterns).
- Another definition of support to better fit the OLAP framework.
- Modular hierarchy management in order to enhance customized and focused knowledge discovery.

| | Agrawal | Jiawei Han | Yu | HYPE |
|-------------------------------|---------|------------|-----|------------|
| Multidimensionality | No | No | Yes | Yes |
| Simulation of multi. | No | ?? | — | Yes |
| Sequential patterns | Yes | No | Yes | Yes |
| Hierarchy in patterns | Several | Single | No | Yes |
| Condensed Representation | No | No | No | No |
| Counting sup. with aggregates | No | No | No | No |
| User interaction | No | No | No | Not enough |

HYPE

Introduction
OLAP & KDD
Sequential Patterns
Multidimensional Framework
Hierarchies
Contributions
Data Model
Definitions
Algorithms
Experiments
Conclusions and Future Work



R. Agrawal and R. Srikant.

Mining sequential patterns.

In *Proc. 1995 Int. Conf. Data Engineering (ICDE'95)*,
pages 3–14, 1995.



J. Han and Y. Fu.

Mining multiple-level association rules in large databases.
IEEE Trans. Knowl. Data Eng., 11(5):798–804, 1999.



M. Plantevit, Y. W. Choong, A. Laurent, D. Laurent, and
M. Teisseire.

M²SP: Mining sequential patterns among several
dimensions.

In *Proc. 2005 PKDD*, pages 205–216, 2005.



C.-C. Yu and Y.-L. Chen.

Mining sequential patterns from multidimensional
sequence data.

IEEE Transactions on Knowledge and Data Engineering,
17(1):136–140, 2005.