

# Pattern-Based Classification of Demographic Sequences

Dmitry I. Ignatov<sup>1</sup>, Danil Gizdatullin<sup>1</sup>, Ekaterina Mitrofanova<sup>1</sup>, Anna Muratova<sup>1</sup>, Jaume Baixeries<sup>2</sup>

<sup>1</sup>National Research University Higher School of Economics, Moscow

<sup>2</sup>Universitat Politècnica de Catalunya, Barcelona

Intelligent Data Processing 2016, Barcelona

# Possible life events

- First job (job)
- The highest education degree is obtained (education)
- Leaving parents' home (separation)
- First partner (partner)
- First marriage (marriage)
- First child birth (children)
- Break-up (parting)
- ... (divorce)

# Data and problem statement

[Ignatov et al., 2015],[Blockeel et al., 2001]

Generation and Gender Survey (GGG): three waves panel data for 11 generations of Russian citizens starting from 30s

## Binary classification

1545 men

3312 women

## Examples of sequential patterns

- $\langle \{ \text{education, separation} \}, \{ \text{work} \}, \{ \text{marriage} \}, \{ \text{children} \} \rangle (m)$
- $\langle \{ \text{work} \}, \{ \text{marriage} \}, \{ \text{children} \} \{ \text{education} \} \rangle (f)$
- $\langle \{ \text{partner} \}, \{ \text{marriage, separation} \}, \{ \text{children} \} \rangle (f)$

# Basic definitions

Textbooks of Han et al., Zaki & Meira, Aggarwal et al., etc

- $s = \langle s_1, \dots, s_k \rangle$  is the **subsequence** of  $s' = \langle s'_1, \dots, s'_k \rangle$  ( $s \preceq s'$ ) if  $k \leq k'$  and there exist  $1 \leq r_1 < r_2 < \dots < r_k \leq k'$  such  $s_j = s'_{r_j}$  for all  $1 \leq j \leq k$ .
- $support(s, D)$  is the **support** of a sequence  $s$  in  $D$ , i.e. the number of sequences in  $D$  such that  $s$  is their subsequence.

$$support(s, D) = |\{s' | s' \in D, s \preceq s'\}|$$

- $s$  is a **frequent closed sequence (sequential pattern)** if there is no  $s'$  such that  $s \prec s'$  and

$$support(s, D) = support(s', D)$$

# Example

Let  $D$  be a set of sequences:

Table: Dataset  $D$ .

$s_1$	$\{a, b, c\}\{a, b\}\{b\}$
$s_2$	$\{a\}\{a, c\}\{a\}$
$s_3$	$\{a, b\}\{b, c\}$

- $I = \{a, b, c\}$  is the set of all items (atomic events)
- $\langle\{a, b\}\{b\}\rangle$  belongs to  $s_1$  and  $s_3$  but it is missing in  $s_2$
- $support_D(\langle\{a, b\}\{b\}\rangle) = 2$
- $\{\langle\{a\}\rangle, \langle\{c\}\rangle, \langle\{a\}\{c\}\rangle, \langle\{a, b\}\{b\}\rangle, \langle\{a, c\}\{a\}\rangle\}$  is the set of closed sequences.

## Growth Rate

$$\text{growth\_rate}_{D' \rightarrow D''}(X) = \begin{cases} \frac{\text{supp}_{D''}(X)}{\text{supp}_{D'}(X)} & \text{if } \text{supp}_{D'}(X) \neq 0 \\ 0 & \text{if } \text{supp}_{D''}(X) = \text{supp}(X) = 0 \\ \infty & \text{if } \text{supp}_{D''}(X) \neq 0 \text{ and } \text{supp}_{D'}(X) = 0 \end{cases}$$

## Class score

$$\text{score}(s, C) = \sum_{e \subseteq s, e \in E(c)} \frac{\text{growth\_rate}_C(e)}{\text{growth\_rate}_C(e) + 1} \cdot \text{supp}_C(e)$$

## Score normalization

$$\text{normal\_score}(s, C) = \frac{\text{score}(s, C)}{\text{median}(\{\text{growth\_rate}_C(e_i)\})}$$

## Classification rule

$$\text{class}(s) = \begin{cases} C_1, & \text{if } \text{normal\_score}(s, C_1) > \text{normal\_score}(s, C_2) \\ C_2, & \text{if } \text{normal\_score}(s, C_1) < \text{normal\_score}(s, C_2) \\ \text{undetermined} & \text{if } \text{normal\_score}(s, C_1) = \text{normal\_score}(s, C_2) \end{cases}$$

- $s = \langle s_1, \dots, s_k \rangle$  is a **gapless prefix-based subsequence** of  $s' = \langle s'_1, \dots, s'_{k'} \rangle$  ( $s * = s'$ ) if  $k \leq k'$  and  $\forall i \in k' : s_i = s'_i$ .
- **Support of gapless prefix-based sequences**  
Let  $T$  be a set of sequences.

$$\text{support}(s, T) = \frac{|\{s' | s' \in T, s * = s'\}|}{|T|}$$

- Let  $0 < \text{minSup} \leq 1$  be a minimal support parameter and  $D$  is a set of sequences then **searching for prefix-based gapless sequential patterns** is the task of enumeration of all prefix-based gapless sequences  $s$  such that  $\text{support}(s, D) \geq \text{minSup}$ . Every sequence  $s$  with  $\text{support}(s, D) \geq \text{minSup}$  is called a **prefix-based gapless sequential pattern**.
- Prefix-based gapless sequential pattern (PGSP)  $p$  is called **closed** if there is no PGSP  $d$  of greater or equal support such that  $d = p*$ .

## Example

Table:  $D$  is a set of sequences.

$s_1$	$\{a\}\{b\}\{d\}$
$s_2$	$\{a\}\{b\}\{c\}$
$s_3$	$\{a, b\}\{b, c\}$

$$s = \langle \{a\}\{b\} \rangle$$

- $I = \{a, b, c\}$  is the set of all items (atomic events)
- $s_1 = s^*$ ;  $s_2 = s^*$
- $s_3 \neq s^*$
- $Supp_D(s) = \frac{2}{3}$
- $\langle \{a\}\{b\} \rangle$  is closed,  $\langle \{a\} \rangle$  is not closed.

- $(S, (D, \sqsubseteq), \delta)$  is a pattern structure
- $S$  is a set of objects,  $D$  is a set of their possible descriptions
- $\delta(g)$  is the description of  $g$  from  $S$
- Galois connection is given by  $\diamond$  operator as follows:

$$A^\diamond := \bigsqcap_{g \in A} \delta(g) \text{ for } A \subseteq S$$

$$d^\diamond := \{s \in S \mid d \sqsubseteq \delta(g)\} \text{ for } d \in D$$

- For two sequences  $\sqsubseteq$  may result in their largest common prefix subsequence

A pair  $(A, d)$  is called a **pattern concept** of a pattern structure  $(S, (D, \sqcap), \delta)$  if

- 1  $A \subseteq S$
- 2  $d \in D$
- 3  $A^\diamond = d$
- 4  $d^\diamond = A$

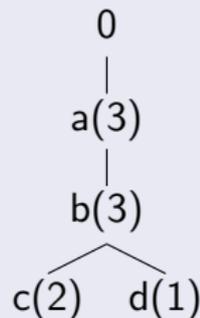
## Example

$s_1 : \langle a, b, c \rangle$

$s_2 : \langle a, b, c \rangle$

$s_3 : \langle a, b, d \rangle$

## Tree



## Pattern concepts (PCs)

$(\{s_1, s_2, s_3\}, \langle a, b \rangle); (\{s_1, s_2\}, \langle a, b, c \rangle)$

$(\{s_1\}, \langle a, b, c \rangle)$  is not a PC;  $(\{s_3\}, \langle a, b, d \rangle)$

# Pattern-based JSM-hypotheses

[Finn, 1981], [Kuznetsov, 1993], [Ganter et al, 2004]

## Positive, negative and undetermined pattern structures

$$\mathbb{K}_{\oplus} = (S_{\oplus}, (D, \sqcap), \delta_{\oplus})$$

$$\mathbb{K}_{\ominus} = (S_{\ominus}, (D, \sqcap), \delta_{\ominus})$$

There is a pattern structure of undetermined examples:

$$\mathbb{K}_{\tau} = (S_{\tau}, (D, \sqcap), \delta_{\tau})$$

## Hypothesis

A **hypothesis** is a pattern intent that belongs to examples from a fixed class only

A pattern intent  $h$  is a positive hypothesis (dually for negative hypotheses) if

$$\forall s \in S_{\ominus} (s \in S_{\oplus}) : h \not\sqsubseteq s^{\ominus} (h \not\sqsubseteq s^{\oplus})$$

# Hypotheses generation: An example

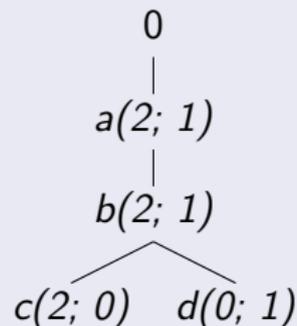
## Sequential classification rules

$s_1 : \langle a, b, c \rangle - \text{class } 0$

$s_2 : \langle a, b, c \rangle - \text{class } 0$

$s_3 : \langle a, b, d \rangle - \text{class } 1$

## Prefix-tree



## Hypotheses

$\langle \{a\}, \{b\}, \{c\} \rangle$  is a hypothesis of class 0

$\langle \{a\}, \{b\}, \{d\} \rangle$  is a hypothesis of class 1

# Classification via hypotheses

$$\text{class}(g_\tau) = \begin{cases} \text{positive} & \text{if } \exists h_\oplus, h_\oplus \sqsubseteq \delta(g_\tau) \text{ and } \nexists h_\ominus, h_\ominus \sqsubseteq \delta(g_\tau) \\ \text{negative} & \text{if } \nexists h_\oplus, h_\oplus \sqsubseteq \delta(g_\tau) \text{ and } \exists h_\ominus, h_\ominus \sqsubseteq \delta(g_\tau) \\ \text{undetermined} & \text{if } \exists h_\oplus, h_\oplus \sqsubseteq \delta(g_\tau) \text{ and } \exists h_\ominus, h_\ominus \sqsubseteq \delta(g_\tau) \\ \text{undetermined} & \text{if } \nexists h_\oplus, h_\oplus \sqsubseteq \delta(g_\tau) \text{ and } \nexists h_\ominus, h_\ominus \sqsubseteq \delta(g_\tau) \end{cases}$$

## Growth Rate

$$\text{GrowthRate}(g, \mathbb{K}_{\oplus}, \mathbb{K}_{\ominus}) = \frac{\text{Sup}_{\mathbb{K}_{\oplus}}(g)}{\text{Sup}_{\mathbb{K}_{\ominus}}(g)}$$

## Emerging patterns

A pattern is called **emerging pattern** if its growth rate is greater than or equal to  $\Theta_{min}$

$$\text{GrowthRate}(g, \mathbb{K}_{\oplus}, \mathbb{K}_{\ominus}) > \Theta_{min}$$

# Emerging patterns for classification

$s$  is a new object

$$\text{normal\_score}_{\oplus}(s) = \frac{\sum_{p \in P_{\oplus}} \text{GrowthRate}(p, \mathbb{K}_{\oplus}, \mathbb{K}_{\ominus})}{\text{median}(\text{GrowthRate}(P_{\oplus}))} : p \sqsubseteq s$$

$$\text{normal\_score}_{\ominus}(s) = \frac{\sum_{p \in P_{\ominus}} \text{GrowthRate}(p, \mathbb{K}_{\ominus}, \mathbb{K}_{\oplus})}{\text{median}(\text{GrowthRate}(P_{\ominus}))} : p \sqsubseteq s$$

Classification via emerging patterns

$$\text{class}(s) = \begin{cases} \text{positive} & \text{if } \text{normal\_score}_{\oplus}(s) > \text{score}_{\ominus}(s) \\ \text{negative} & \text{if } \text{normal\_score}_{\oplus}(s) < \text{score}_{\ominus}(s) \\ \text{undetermined} & \text{if } \text{normal\_score}_{\oplus}(s) = \text{normal\_score}_{\ominus}(s) \end{cases}$$

# Classification algorithm for gapless prefix-based sequential patterns

- 1 Build the prefix tree for the input sequences.
- 2 For each tree node calculate its Growth Rate.
- 3 For every new sequence traverse the tree and compute the Score for each class.
- 4 Compare the Score value for different classes and classify the new sequence.

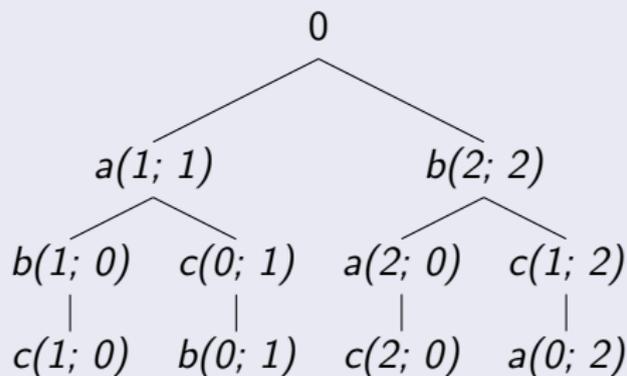
# Execution example

## Input sequences

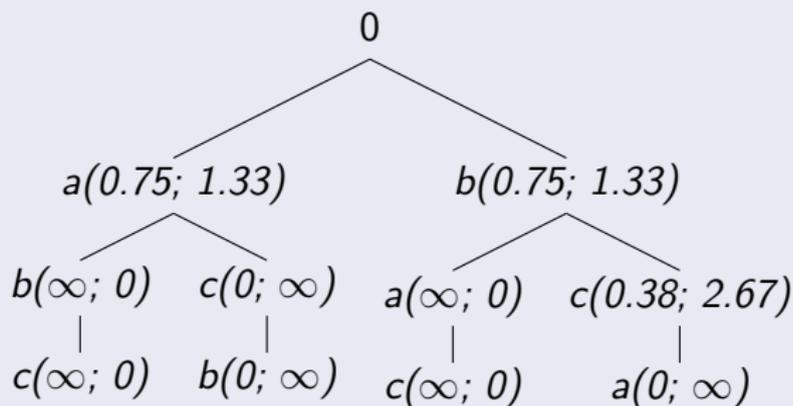
*class\_0* : {⟨{a}{b}{c}⟩, ⟨{b}{a}{c}⟩, ⟨{b}{a}{c}⟩, ⟨{b}{c}⟩}

*class\_1* : {⟨{a}{c}{b}⟩, ⟨{b}{c}{a}⟩, ⟨{b}{c}{a}⟩}

## Prefix tree



## Counting Growth Rate



## New sequence

$\langle \{b\}; \{c\}; \{a\} \rangle - ???$

$$\text{Score}_0 = 0$$

$$\text{Score}_1 = 2.67 + \infty = \infty$$

# Comparison of closed and non-closed patterns

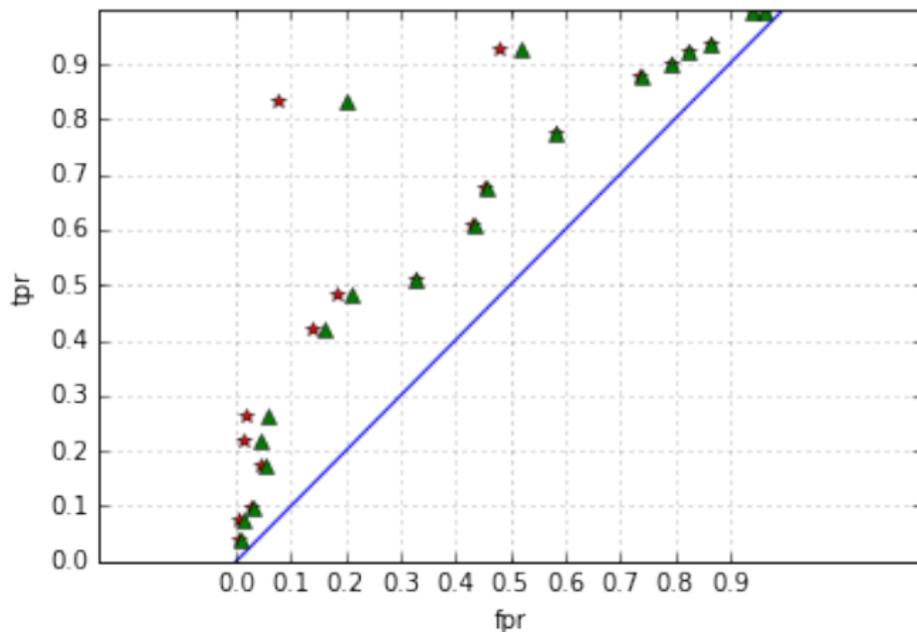


Figure: TPR vs FPR for closed and non-closed patterns

# Experiments and results

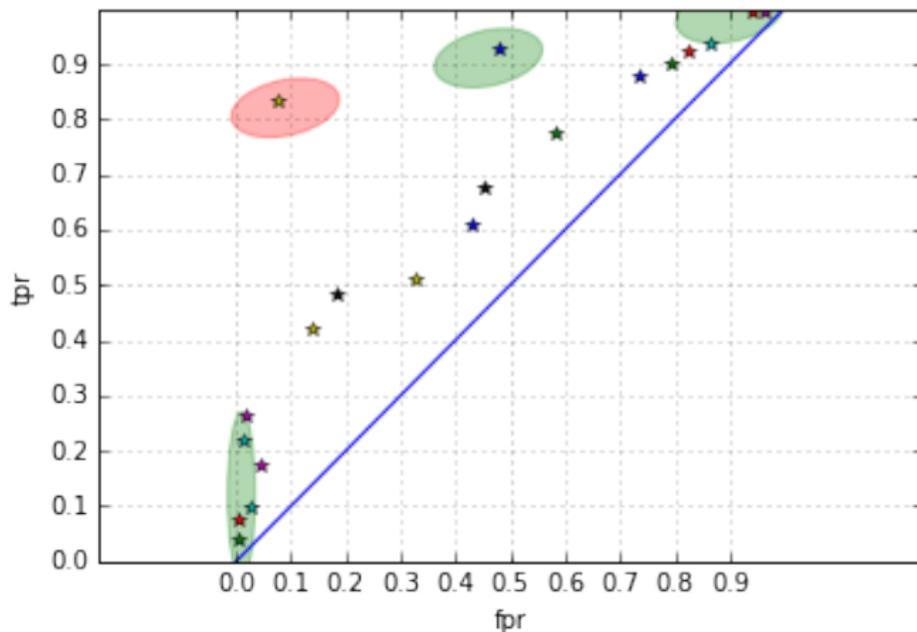


Figure: TPR-FPR for classification via gapless prefix-based patterns

# Interesting patterns (women)

$(\langle\{work, separation\}, \{marriage\}, \{children\}, \{education\}\rangle, [inf, 0.006])$

$(\langle\{separation, partner\}, \{marriage\}\rangle, [inf, 0.006])$

$(\langle\{work, separation\}, \{marriage\}, \{children\}\rangle, [inf, 0.008])$

$(\langle\{work, separation\}, \{marriage\}\rangle, [inf, 0.009])$

# Interesting patterns (men)

$(\langle\{\textit{education}\}, \{\textit{marriage}\}, \{\textit{work}\}, \{\textit{children}\}, \{\textit{separation}\}\rangle, [10.6, 0.006])$

$(\langle\{\textit{education}\}, \{\textit{marriage}\}, \{\textit{work}\}, \{\textit{children}\}\rangle, [12.7, 0.007])$

$(\langle\{\textit{educ}\}, \{\textit{work}\}, \{\textit{part}\}, \{\textit{mar}\}, \{\textit{sep}\}, \{\textit{ch}\}\rangle, [10.6, 0.006])$

- ① We have studied several pattern mining techniques for demographic sequences including pattern-based classification in particular.
- ② We have fitted existing approaches for sequence mining of a special type (gapless and prefix-based ones).
- ③ The results for different demographic groups (classes) have been obtained and interpreted.
- ④ In particular, a classifier based on emerging sequences and pattern structures has been proposed.

Thank you!

Questions?