

Pattern mining in personal demographic trajectories

Danil Gizdatullin¹, Dmitry I. Ignatov¹

¹National Research University Higher School of Economics, Moscow

27 октября 2017 г.

- 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 (GGS): three waves panel data for 11 generations of Russian citizens starting from 30s

Binary classification

1545 men

3312 women

Examples of sequential patterns

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

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:

Таблица: 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

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

Class score

$$score(s, C) = \sum_{e \subseteq s, e \in E(c)} \frac{growth_rate_C(e)}{growth_rate_C(e) + 1} \cdot supp_C(e)$$

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 $s' = \langle s'_1, \dots, s'_{k'} \rangle$ ($s \leq s'$) 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.

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

- Let $0 < minSup \leq 1$ 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 $support(s, D) \geq minSup$. Every sequence s with $support(s, D) \geq minSup$ 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

Таблица: 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.

- $(S, (D, \sqcap), \delta)$ 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 \sqcap 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 $d^\diamond = A$

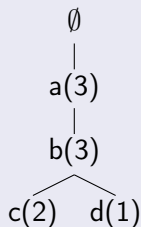
Example

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

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

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

Tree



Pattern concepts (PCs)

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

$(\{s_3\}, \langle a, b, d \rangle)$

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

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

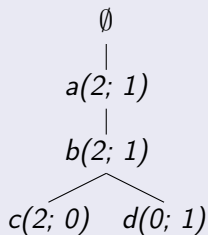
Sequential classification rules

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

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

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

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

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

Emerging patterns

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

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

Emerging patterns for classification

s is a new object

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

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

Classification via emerging patterns

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

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.

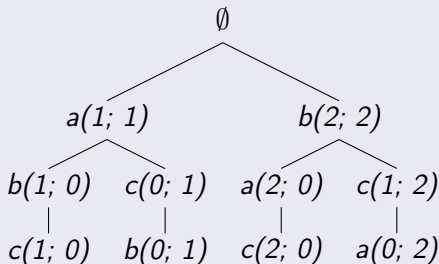
Execution example

Input sequences

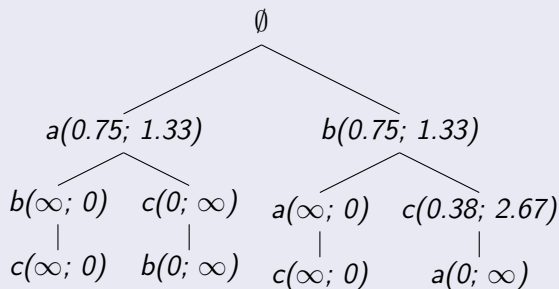
class_0 : {⟨{a}{b}{c}⟩, ⟨{b}{a}{c}⟩, ⟨{b}{a}{c}⟩, ⟨{b}{c}⟩}

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

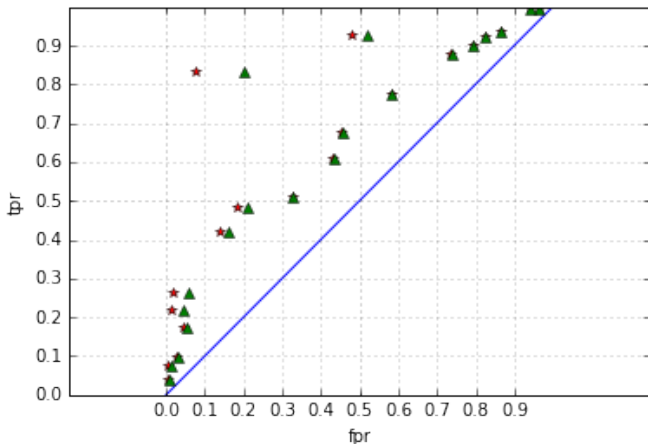


Рис.: TPR vs FPR for closed and non-closed patterns

Experiments and results

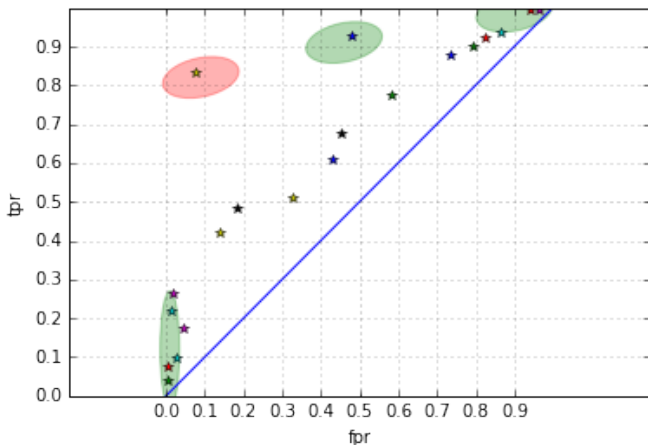


Рис.: TPR-FPR for classification by gender via gapless prefix-based patterns

Interesting patterns (women)

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

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

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

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

Interesting patterns (men)

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

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

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

Experiments and results

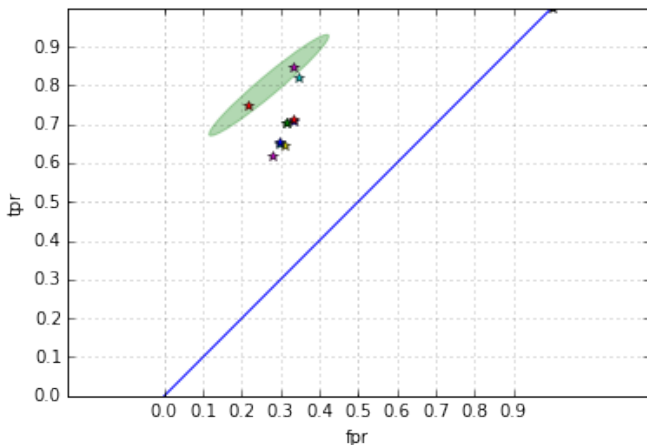


Рис.: TPR-FPR for classification by generation via gapless prefix-based patterns

Interesting patterns (Different Generations; Women)

Old women

$(\langle\{work\}, \{separation\}\rangle, [1.85, 0.38])$

$(\langle\{work\}, \{marriage, separation\}\rangle, [3.92, 0.08])$

Young women

$(\langle\{education\}\rangle, [1.84, 0.26])$

$(\langle\{education\}, \{work\}\rangle, [4.01, 0.1])$

Experiments and results

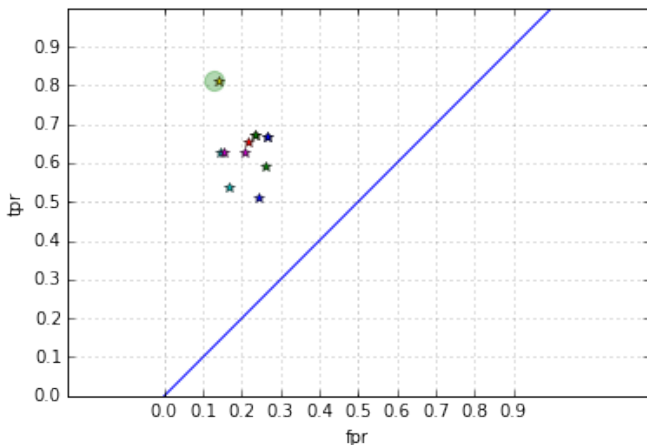


Рис.: TPR-FPR for classification by generation via gapless prefix-based patterns

Interesting patterns (Different Generations; Men)

Old men

$(\langle\{\textit{work}\}, \{\textit{marriage, separation}\}, \{\textit{education}\}\rangle, [13.52, 0.025])$

$(\langle\{\textit{work}\}, \{\textit{marriage}\}, \{\textit{separation}\}\rangle, [22.87, 0.042])$

$(\langle\{\textit{work}\}, \{\textit{marriage}\}, \{\textit{separation}\}, \{\textit{education}\}\rangle, [\infty, 0.0208])$

Young men

$(\langle\{\textit{education}\}, \{\textit{work}\}, \{\textit{separation}\}, \{\textit{marriage}\}, \{\textit{children}\}\rangle, [10.58, 0.020])$

$(\langle\{\textit{education}\}, \{\textit{work}\}, \{\textit{separation, partner}\}, \{\textit{marriage}\}\rangle, [8.65, 0.016])$

$(\langle\{\textit{education}\}, \{\textit{marriage, separation}\}\rangle, [7.69, 0.015])$

- 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?