

# Выделение типов именованных сущностей используя контекстную встречаемость слов

Хайруллин Ринат

Московский физико-технический институт  
Факультет управления и прикладной математики  
Кафедра интеллектуальных систем

Научный руководитель д.ф.-м.н. В.А. Серебряков

Москва,  
2018 г.

## Дано:

$\mathcal{D} = \{d_1, \dots, d_D\}$  – копус текстов,

$\mathcal{V} = \{w_1, \dots, w_N\}$  – представление копуса  $\mathcal{D}$  в виде последовательности слов,

$\mathcal{T} = \{t_1, \dots, t_T\}$  – множество типов именованных сущностей.

## Задача:

Определить тип  $t \in \mathcal{T}$  для каждого слова  $w \in \mathcal{V}$ .

## Проблемы существующих алгоритмов:

- отбор признаков важнее выбора модели,
- неустойчивость к смене предметной области,
- требуется большой объем обучающей выборки.

## Задача

*Построить алгоритм выделения именованных сущностей, обладающий следующими свойствами:*

- 1 *независимость от предметной области,*
- 2 *использование универсальных признаков,*
- 3 *использование частично размеченной выборки.*

## Способ решения

- **независимость от предметной области:**  
*выделение значимых слов и фраз в тексте, которые могут быть именованными сущностями,*
- **универсальные признаки:**
  - 1 *выделение слов и фраз, которые задают отношения между словами и фразами в предложении,*
  - 2 *кластеризация выделенных отношений на основе контекста,*
- **использование частичной разметки выборки**

**Дано:**  $\mathcal{D}$ ,  $(\Psi, \mathcal{T}_\Psi, \sigma)$ ,  $\mathcal{T} \subset \mathcal{T}_\Psi$

$\phi = \{w_{d,i}, \dots, w_{d,i+k}\}$   $k \geq 0$  – значимая фраза, если  $\rho(\phi) > \alpha$   
 $\rho(\cdot)$  – некоторая функция значимости,  $\alpha$  – порог значимости

**Задача:**

- построить множество значимых фраз  $\mathcal{Q}$

$$\mathcal{M}_d(q) = \{i \mid \phi_i = q\} \quad i \in \{\tilde{n}_{d-1} + 1, \dots, \tilde{n}_d\}$$

$$\mathcal{M} = \bigcup_{d \in \mathcal{D}} \bigcup_{q \in \mathcal{Q}} \mathcal{M}_d(q)$$

- построить множество отношений  $\mathcal{R}$

$$\mathcal{L}_d(r) = \{i \mid \phi_i = r\} \quad i \in \{\tilde{n}_{d-1} + 1, \dots, \tilde{n}_d\}$$

$$\mathcal{L} = \bigcup_{d \in \mathcal{D}} \bigcup_{r \in \mathcal{R}} \mathcal{L}_d(r)$$

- выделить подмножество  $\mathcal{Q}_L = \{q \in \mathcal{Q} \mid q \in \Psi, \mathcal{T}(q) \cap \mathcal{T} = \{t\}\}$
- для всех  $q \in \mathcal{Q}_U = \mathcal{Q} \setminus \mathcal{Q}_L$  в корпусе  $\tilde{\mathcal{D}} = \text{sort}(\mathcal{M} \cup \mathcal{L})$

$$t(q_i) = \underset{j \in \mathcal{T}}{\operatorname{argmax}} \mathbf{y}(q_i)_j \quad \mathbf{y}(q_i) \in [0, 1]^{|\mathcal{T}|} \quad i \in \mathcal{M}(q) \quad j \in \{1, \dots, |\mathcal{T}|\}$$

# Пример графового представления

... [в]:PR [РФ]:EM [Воробьев]:EM [был назначен]:RP ...  
... [Халиков]:EM [был назначен на]:RP должность торгпреда [РФ]:EM  
[в]:RP [Пакистане]:EM ...  
... [Арам Габрелянов]:EM [в]:RP феврале [был назначен]:RP  
[председателем совета директоров]:EM ...

... [пишет в]:RP среду [газета Известия]:EM ...  
... [в]:RP среду [газета Известия]:EM [написала]:RP что ....  
... указанная информация впервые [была распространена]:RP [газетой  
Известия]:EM ...

# Представление данных в виде графа

$$|\mathcal{M}| = m \quad |\mathcal{Q}| = n \quad |\mathcal{R}| = k$$

$\mathcal{G}_Q = [\mathbb{1}_{\{\phi_i=q_j\}}] \in \mathbb{R}^{m \times n}$  – матрица соответствий  $\mathcal{M}$  и  $\mathcal{Q}$

$\mathcal{G}_Z = [\mathbb{1}_{\{\phi_{j+1}=q_i, i,j \in s\}}] \in \mathbb{R}^{m \times k}$  – матрица соответствий слева/справа  $\mathcal{M}$  и  $\mathcal{R}$  в предложении  $s$ , где  $Z \in \{L, R\}$

$$\mathcal{W}_Z = \mathcal{G}_Q^T \mathcal{G}_Z \in \mathbb{R}^{n \times k}$$

$\mathcal{W}_M \in \mathbb{R}^{n \times n}$  – матрица соответствий между значимыми фразами

Для каждого  $\phi_i \in \mathcal{M}$  определим вектор весов  $q_j$ -х, которые встречаются с  $\phi_i$  в одном предложении  $s$

$$\mathbf{f}_{\phi_i} = \{\nu_s(q_j) \log\left(\frac{|\mathcal{D}|}{\nu_{\mathcal{D}}(q_j)}\right) \mid \phi_i, q_j \in s\} \in \mathbb{R}^n$$

$$W_{\mathcal{M},i,j} = \begin{cases} \exp\left(\frac{-\|\mathbf{f}_{\phi_i} - \mathbf{f}_{\phi_j}\|}{t}\right) & \text{если } \mathbf{f}^{(i)} \in N_k(\mathbf{f}^{(j)}) \text{ или } \mathbf{f}^{(j)} \in N_k(\mathbf{f}^{(i)}) \\ & \text{и } \phi_i = \phi_j \\ 0 & \text{иначе} \end{cases}$$

## Hypothesis (H1)

*Если фраза  $\phi$  часто встречается слева/справа от отношения  $r$  в предложениях, то  $\phi$  и  $r$  должны относиться к одному типу  $t$ .*

## Hypothesis (H2)

*Если в контекстах двух фраз  $\phi_i$  и  $\phi_j$  много общих слов и  $\phi_i = \phi_j$ , то  $\phi_i$  и  $\phi_j$  должны относиться к одному типу  $t$ .*

## Hypothesis (H3)

*Два отношения  $r_i$  и  $r_j$  относятся к одному кластеру если у них много общих контекстных слов.*

## Hypothesis (H4)

*Два отношения  $r_i$  и  $r_j$  относятся к одному типу  $t$ , если они принадлежат одному кластеру.*

$\phi_i \in \mathcal{M}$ :  $\mathbf{y} \in \mathbb{R}^{m \times T}$ ,  $q_i \in \mathcal{Q}$ :  $\mathbf{c} \in \mathbb{R}^{n \times T}$ ,  $r_i \in \mathcal{R}$ :  $\mathcal{P}_Z \in \mathbb{R}^{k \times T}$

$\alpha, \gamma, \mu \in [0, 1]$

$$\mathcal{F}(\mathbf{c}, \mathcal{P}_L, \mathcal{P}_R) = \sum_{Z \in \{L, R\}} \sum_{i=1}^n \sum_{j=1}^k W_{Z,i,j} \left\| \frac{\mathbf{c}_i}{\sqrt{\mathbf{D}_{Q,Z,i,i}}} - \frac{\mathcal{P}_{Z,j}}{\sqrt{\mathbf{D}_{R,Z,j,j}}} \right\|_2^2 \quad (1)$$

$$\begin{aligned} \Omega_{\gamma, \mu}(\mathbf{y}, \mathbf{c}, \mathcal{P}_L, \mathcal{P}_R) = & \frac{1}{2} \|\mathbf{y} - f(\mathcal{G}_Q \mathbf{c}, \mathcal{G}_L \mathcal{P}_L, \mathcal{G}_R \mathcal{P}_R)\|_F^2 + \\ & \frac{\gamma}{2} \sum_{i,j=1}^n W_{\mathcal{M},i,j} \left\| \frac{\mathbf{y}_i}{\sqrt{\mathbf{D}_{\mathcal{M},i,i}}} - \frac{\mathbf{y}_j}{\sqrt{\mathbf{D}_{\mathcal{M},j,j}}} \right\|_2^2 + \frac{\mu}{2} \|\mathbf{y} - \mathbf{y}_0\|_F^2 \end{aligned} \quad (2)$$

$$\mathcal{L}_\alpha(\mathcal{P}_L, \mathcal{P}_R, \{\mathbf{U}_v, \mathbf{V}_v, \beta_v\}, \mathbf{V}^*) = \sum_{v=0}^d \beta_v (\|\mathbf{F}_v^T - \mathbf{U}_v \mathbf{V}_v^T\|_F^2 + \alpha \|\mathbf{V}_v \mathbf{Q}_v - \mathbf{V}^*\|_F^2) \quad (3)$$

$$\mathbf{D}_{Q,Z} = \sum_{i=1}^n W_{Z,i,j}, \quad \mathbf{D}_{R,Z} = \sum_{j=1}^k W_{Z,i,j}, \quad \mathbf{S}_Z = \mathbf{D}_{Q,Z}^{-\frac{1}{2}} \mathbf{W}_Z \mathbf{D}_{R,Z}^{-\frac{1}{2}}$$

Аналогично определяется  $\mathbf{S}_{\mathcal{M}}$



$$\begin{aligned}
 & \min_{\mathcal{Y}, \mathcal{C}, \mathcal{P}_L, \mathcal{P}_R, \{\mathbf{U}_v, \mathbf{V}_v, \beta_v\}, \mathbf{V}^*} (\mathcal{F} + \Omega_{\gamma, \mu} + \mathcal{L}_\alpha) \\
 & \text{s.t. } \{\mathbf{U}_v, \mathbf{V}_v\}, \mathbf{V}^* \geq 0, \quad \sum_{v=1}^d \exp(-\beta_v) = 1
 \end{aligned} \tag{4}$$

Решение: Определим тип  $t(q_i) = \operatorname{argmax} \mathcal{Y}_i$

Оценки параметров:

$$\mathcal{Y} = [(1 + \gamma + \mu)\mathbf{1} - \gamma \mathcal{S}_M]^{-1} [f(\mathcal{G}_Q \mathcal{C}, \mathcal{G}_L \mathcal{P}_L, \mathcal{G}_R \mathcal{P}_R)^q + \mu \mathcal{Y}_0^q] \tag{5}$$

$$\mathcal{C} = \frac{1}{2} [\mathcal{S}_L \mathcal{P}_L + \mathcal{S}_R \mathcal{P}_R + \mathcal{G}_Q^T (\mathcal{Y} - \mathcal{G}_L \mathcal{P}_L - \mathcal{G}_R \mathcal{P}_R)] \tag{6}$$

$$\mathcal{P}_L = [(1 + \beta_0)\mathbf{1} + \mathcal{G}_L^T \mathcal{G}_L]^{-1} [\mathcal{S}_L^T \mathcal{C} + \mathcal{G}_L^T [\mathcal{Y} - \mathcal{G}_Q \mathcal{C} - \mathcal{G}_R \mathcal{P}_R] + \beta_0 \mathbf{V}_0 \mathbf{U}_0^T] \tag{7}$$

$$\mathcal{P}_R = [(1 + \beta_1)\mathbf{1} + \mathcal{G}_R^T \mathcal{G}_R]^{-1} [\mathcal{S}_R^T \mathcal{C} + \mathcal{G}_R^T [\mathcal{Y} - \mathcal{G}_Q \mathcal{C} - \mathcal{G}_L \mathcal{P}_L] + \beta_1 \mathbf{V}_1 \mathbf{U}_1^T] \tag{8}$$

# Кластеризация на множестве отношений

$|\mathcal{R}| = k, r_i \in \mathcal{R}: \mathcal{P}_Z \in \mathbb{R}^{k \times T}, \mathbf{F}_{N_s} \in \mathbb{R}^{k \times n_s}, \mathbf{F}_{N_c} \in \mathbb{R}^{k \times n_c}$   
 $\{\mathbf{F}_0, \mathbf{F}_1, \mathbf{F}_2, \mathbf{F}_3\} = \{\mathcal{P}_L, \mathcal{P}_R, \mathbf{F}_{N_s}, \mathbf{F}_{N_c}\}$  – признаки представления  
 $\mathbf{U}_v \in \mathbb{R}_{\geq 0}^{N_v \times C}, \mathbf{V}_v \in \mathbb{R}_{\geq 0}^{k \times C}, \mathbf{V}^* \in \mathbb{R}_{\geq 0}^{k \times C} \quad \mathbf{Q} \in \mathbb{R}_{\geq 0}^{C \times C} \quad Q_{v,c,c} = \sum_{i=1}^{N_v} U_{v,i,c}$

$$\min_{\{\mathbf{U}_v, \mathbf{V}_v, \beta_v\}, \mathbf{V}^*} \sum_{v=0}^d \beta_v \underbrace{(\|\mathbf{F}_v^T - \mathbf{U}_v \mathbf{V}_v^T\|_F^2 + \alpha \|\mathbf{V}_v \mathbf{Q}_v - \mathbf{V}^*\|_F^2)}_{\delta_v}$$

$$\text{s.t. } \{\mathbf{U}_v, \mathbf{V}_v\}, \mathbf{V}^* \geq 0 \quad \sum_{v=1}^d \exp(-\beta_v) = 1$$

$$U_{i,k} \leftarrow U_{i,k} \frac{(\mathbf{FV})_{i,k} + \alpha \sum_{j=1}^k V_{j,k} V_{j,k}^*}{(\mathbf{UV}^T \mathbf{V})_{i,k} + \alpha \sum_{l=1}^{N_v} U_{l,k} \sum_{j=1}^k V_{j,k}^2} \quad (9)$$

$$\mathbf{U} \leftarrow \mathbf{UQ}^{-1} \quad \mathbf{V} \leftarrow \mathbf{VQ} \quad (10)$$

$$V_{i,k} \leftarrow V_{i,k} \frac{(\mathbf{F}^T \mathbf{U})_{i,k} + \alpha V_{j,k}^*}{(\mathbf{VU}^T \mathbf{U})_{j,k} + \alpha V_{j,k}} \quad (11)$$

$$\mathbf{V}^* = \frac{\sum_{v=1}^d \alpha \mathbf{V}_v \mathbf{Q}_v}{d\alpha} \quad \beta_v = -\log \frac{\delta_v}{\sum_{v=1}^d \delta_v} \quad (12)$$

$$\mathbf{U}_0 \mathbf{Q}_0^{-1} \approx \mathbf{U}_1 \mathbf{Q}_1^{-1} \approx [p(t_i | c)], \quad \mathbf{Q}_0 \mathbf{V}_0^T \approx \mathbf{Q}_1 \mathbf{V}_1^T \approx [p(r_i, c)]$$

---

**Algorithm 1** Optimization procedure

---

- 1:  $\{\mathcal{Y}, \mathcal{C}, \mathcal{P}_L, \mathcal{P}_R\} \leftarrow \{\mathcal{Y}_0, \mathcal{G}_Q^T \mathcal{Y}_0, \mathcal{G}_L^T \mathcal{Y}_0, \mathcal{G}_R^T \mathcal{Y}_0\}$
  - 2:  $\{\mathbf{U}_v, \mathbf{V}_v, \beta_v\}, \mathbf{V}^* \leftarrow$  положительные числа
  - 3: **repeat**
  - 4:    $\mathcal{Y}, \mathcal{C}, \mathcal{P}_L, \mathcal{P}_R \leftarrow Eq.(5 - 8)$
  - 5:   **for**  $v = 0$  to 3 **do**
  - 6:     **repeat**
  - 7:        $\mathbf{U} \leftarrow Eq.(9)$
  - 8:        $\mathbf{U}, \mathbf{V} \leftarrow Eq.(10)$
  - 9:        $\mathbf{V} \leftarrow Eq.(11)$
  - 10:     **until**  $\delta_v$  converges
  - 11:    **endfor**
  - 11:     $\mathbf{V}^*, \beta_v \leftarrow Eq.(12)$
  - 12: **until** Objective in Eq.(4) converges
-

## Признаки:

$\mathcal{N}_s = \{w_1, \dots, w_{n_s}\}$  – все уникальные слова входящие в  $\mathcal{R}$

$\mathcal{N}_c = \{\tilde{w}_1, \dots, \tilde{w}_{n_c}\}$  – все уникальные слова в окне из 10 слов для фраз отношений

Определим матрицы  $\mathbf{F}_{\mathcal{N}_s} \in \mathbb{R}^{k \times n_s}$ ,  $\mathbf{F}_{\mathcal{N}_c} \in \mathbb{R}^{k \times n_c}$ , как:

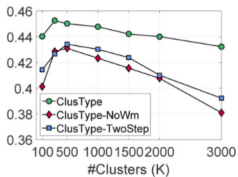
$$F_{\mathcal{N}_s, i, j} = \begin{cases} 1 & \text{если } w_j \in r_i \\ 0 & \text{иначе} \end{cases}$$

$$F_{\mathcal{N}_c, i, j} = TFIDF(\tilde{w}_{i, j})$$

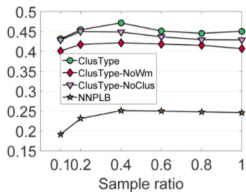
$|\mathcal{V}| = 300000$ ,  $|\mathcal{Q}| = 10000$ ,  $|\mathcal{R}| = 6000$ ,  $|\mathcal{M}| = 31000$ .  
 $|\mathcal{M}_0| = 2000$

Table: Результаты на корпусе FactRuEval и LABINFORM

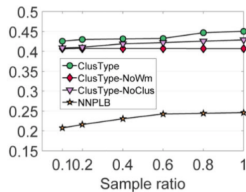
	Type	precision	recall	F1
No Clus	Pers	0.39	0.54	0.45
No Clus	Loc	0.25	0.16	0.2
No Clus	Org	0.32	0.3	0.31
Clus	Pers	0.38	0.65	0.48
Clus	Loc	0.25	0.1	0.15
Clus	Org	0.32	0.24	0.27



(a) #Clusters



(b) Seed set size



(c) Corpus size