

# Text mining

## Извлечение информации и классификация последовательностей

Екатерина Черняк

[echernyak@hse.ru](mailto:echernyak@hse.ru)

Национальный Исследовательский Университет – Высшая Школа Экономики  
НУЛ Интеллектуальных систем и структурного анализа

December 15, 2017

- 1 Извлечение информации
- 2 Извлечение информации по словарям, регулярными выражениями и по шаблонам
- 3 Классификация чанков [chunk]
- 4 Условные случайные поля
- 5 Задача классификации последовательности
  - Условные случайные поля
  - Рекуррентные нейронные сети
  - Заполнение слотов

# Извлечение информации [Information extraction, IE]

Извлечение значимых элементов текста:

- Даты, номера телефонов, адреса
- Именованные сущности [Named entity recognition, NER]
- Отношения, факты и события
- Термины
- Разрешение кореференции [Coreference resolution]
- OpenIE (не сейчас)

# Методы извлечения информации

- Словари, регулярные выражения и шаблоны
- Классификация чанков [chunk]
- Классификация последовательностей [sequence labelling]
- Заполнение слотов [slot filling]

- Wikipedia / DBPedia и другие Wiki
- Словари имен, терминов
- API карт (например, API KudaGo)
- Государственные реестры

# Регулярные выражения

## Номер телефона

```
re.compile('[+0-9 \- \(\ )]8,')
```

## email

```
r"^[a-zA-Z0-9_\.+-]+[a-zA-Z0-9-]+\.[a-zA-Z0-9-]+\.$"
```

## Римские числа

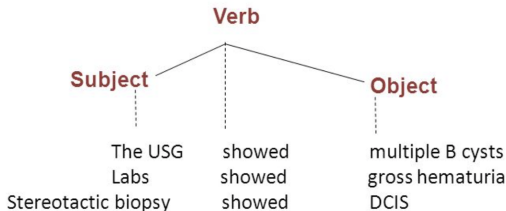
```
^ M{0,4}(CM|CD|D?C0,3)(XC|XL|L?X{0,3})(IX|IV|V?I{0,3})$
```

Шаблоны М. Hearst [Hearst, 1992]:

- ① NP such as NP,\* (and|or) NP
  - ② such NP as NP,\* (and|or) NP
  - ③ NP ,NP\* , (and|or) other NP
  - ④ NP including NP,\* NP (and|or) NP
  - ⑤ NP especially NP,\* (and|or) NP
- Such injuries as bruises, wounds and broken bones
  - Мало обратимые длительные психические расстройства, такие как нарушения или задержка умственного развития

Легко описать с помощью Томита-парсера.

## Medical Semantics: Triples



Источник: <http://slideplayer.com/slide/11723219/>  
Легко описать в грамматике универсальных зависимостей.



# Классификация чанков [chunk]

Чанк [chunk] – последовательность из нескольких токенов, иногда – синтаксическая группа

Именная группа [noun phrase, NP] – Adj\* N\*

Глагольная группа [verb phrase, VP] – Adv\* V

Классификация чанков: данный чанк –

- название города?
- имя человека?
- название компании?
- название группы?

Можно использовать любой бинарный классификатор. Признаки:

- Морфологические тэги
- Регистр
- Есть ли одноименная статья в Википедии
- Слова в левом и правом окнах

# Извлечение отношений [relation extraction]

- Родился-В: **кто** родился **где**?
- Супруги: **кто** женат / замужем **на ком** / **за кем**?
- CEO: **кто** руководит какой **компанией**?
- Сделка: **кто** заключил сделку **с кем**?
- Образование: **образование** обнаружено в **каком органе**?
- Побочный эффект: **препарат** вызвал **побочный эффект**?
- Мутация: **мутация** найдена в **каком гене**?

Задача классификации: дано два слова / чанка,  $c_1, c_2$ .

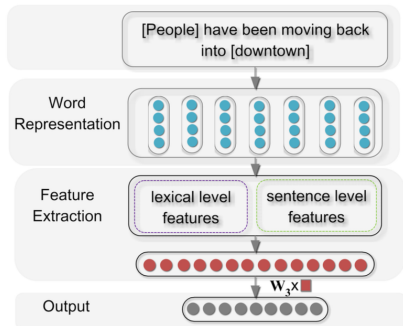
- простая задача бинарной классификации: находятся ли они в отношении  $R(c_1, c_2)$ ?
- задача классификации на несколько классов: в каком отношении  $R$  находятся  $(c_1, c_2)$ ?

The most common audits were about waste and recycling .	Message-Topic(e1,e2)	3	6
The company fabricates plastic chairs .	Product-Producer(e2,e1)	1	4
The school master teaches the lesson with a stick .	Instrument-Agency(e2,e1)	2	8

# Методы классификации в задаче извлечения отношений

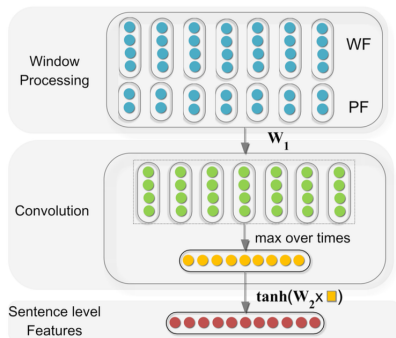
- Любой метод классификации и множество лингвистических признаков: POS-тэги, длина синтаксического пути между чанками, расстояние в WordNet и др.
- CNN и positional embeddings [Zeng et al., 2014, Relation Classification via Convolutional Deep Neural Network]
- Att-BiLSTM и index embeddings [Zhou et al., 2016, Attention-Based Bidirectional Long Short-Term Memory Networks for Relation Classification]

# CNN и positional embeddings [Zeng et al., 2014]



- Вход: слова и word embeddings
- Извлечение признаков: признаки на уровне слова и на уровне предложения
- Выход: число нейронов = количество различных отношение

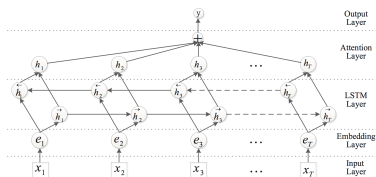
# CNN и positional embeddings [Zeng et al., 2014]



Position embeddings состоят из двух частей: расстояние от текущего слова до  $e_1$  и до  $e_2$ . Каждому числу в соответствие ставится свой вектор. Итоговый вектор слова состоит из трех частей: эмбединг слова, вектор, кодирующий расстояние до  $e_1$  и вектор, кодирующий расстояние до  $e_2$ .



# Att-BiLSTM и index embeddings [Zhou et al., 2016]



Аналогичная идея: вектор слова состоит из эмбединга слова и one-hot кодирования его индекса в словаре.



# Извлечение фактов и событий [facts and events extraction]

Christopher Manning



## Landscape of IE Tasks: Arity of relation

Jack Welch will retire as CEO of General Electric tomorrow. The top role at the Connecticut company will be filled by Jeffrey Immelt.

### Single entity

**Person:** Jack Welch

**Person:** Jeffrey Immelt

**Location:** Connecticut

### Binary relationship

**Relation:** Person-Title

**Person:** Jack Welch

**Title:** CEO

**Relation:** Company-Location

**Company:** General Electric

**Location:** Connecticut

### N-ary record

**Relation:** Succession

**Company:** General Electric

**Title:** CEO

**Out:** Jack Welch

**In:** Jeffrey Immelt

*"Named entity" extraction*

Slide by Andrew McCallum. Used with permission.

Slide by Christopher Manning. Used without permission.

Факты и события:  $N$  – арные отношения

- Automatic Content Extraction (ACE)

- ▶ Тип события (“смерть”)
- ▶ Аргументы события (кто, где, когда, причина)
- ▶ Триггеры (умер, погиб, убит, смерть, убийство)
- ▶ Атрибуты события (модальность, поляризация, время, и др).

- FactRuEval

- ▶ Тип факта (“должность”)
- ▶ Аргументы факта (кто, должность, компания)

- 1 Извлечение информации
- 2 Извлечение информации по словарям, регулярными выражениями и по шаблонам
- 3 Классификация чанков [chunk]
- 4 Условные случайные поля
- 5 Задача классификации последовательности
  - Условные случайные поля
  - Рекуррентные нейронные сети
  - Заполнение слотов

# Задача классификации последовательности

	Британская	актриса	и	крестница	принца	Чарльза	Тара	Томкинсон
POS	Прил.	Сущ.	Союз	Сущ.	Сущ.	Им.Собств.	Им. Собств.	Им. Собств.
IOB (NE)	O	O	O	O	O	B-Per	B-Per	I-Per
IOBES (NE)	O	O	O	O	O	S-Per	B-Per	E-Per
IOBES (R)	O	O	O	O	O	B-Per-1	B-Per-2	I-Per-2
	была	найдена	мертвой	в	ее	квартире	в	Лондоне
POS	Глаг.	Кр. Прич.	Прил.	Пред.	Мест.	Сущ.	Пред.	Им. Собств.
IOB (NE)	O	O	O	O	O	O	O	B-Loc
IOBES (NE)	O	O	O	O	O	O	O	S-Loc
IOBES (R)	O	O	O	O	O	O	O	O
	,	сообщает	BBC	.				
POS	Пункт.	Глаг.	Им. Собств	Пункт.				
IOB (NE)	O	O	B-Org	O				
IOBES (NE)	O	O	S-Org	O				
IOBES (R)	O	O	O	O				

# Определение

Обучающие данные:

- $\mathbf{x} = x_1, x_2, \dots, x_n, x_i \in V, V$  – словарь
- $\mathbf{y} = y_1, y_2, \dots, y_n, y_i \in \{1, \dots, L\}$  – метки
- $\{(\mathbf{x}^{(1)}, \mathbf{y}^{(1)}), (\mathbf{x}^{(2)}, \mathbf{y}^{(2)}), \dots, (\mathbf{x}^{(m)}, \mathbf{y}^{(m)})\}$  – обучающие данные
- экспоненциальная сложность: если длина входной последовательности  $= n$ , всего возможно  $L^n$  решений

Требуется обучить классификатор:  $\mathbf{x} \rightarrow \mathbf{y}$

- $y$  – последовательность
- $y$  – дерево (парсинг)

- Sequence labelling

- ▶ Марковские модели максимальной энтропии [Maximum-entropy Markov model, MEMM]
- ▶ Условные случайные поля [Conditional random fields, CRF]
- ▶ Рекуррентные нейронные сети (biLSTM)
- ▶ (CNN-)biLSTM-CRF [Ma and Hovy, 2016 End-to-end Sequence Labeling via Bi-directional LSTM-CNNs-CRF]

- Structured prediction

- ▶ SVM<sup>struct</sup>
- ▶ Structured perceptron

- Slot filling

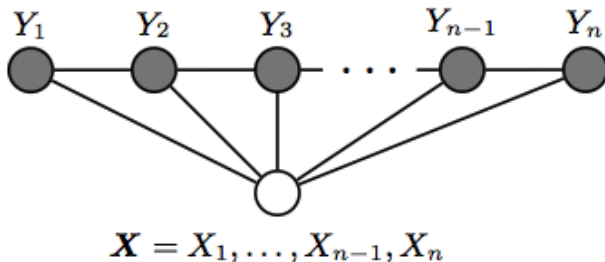
- ▶ biLSTM-CNN-CRF with attention [Liu and Lane, 2016 Attention-Based Recurrent Neural Network Models for Joint Intent Detection and Slot Filling]

- 1 Извлечение информации
- 2 Извлечение информации по словарям, регулярными выражениями и по шаблонам
- 3 Классификация чанков [chunk]
- 4 Условные случайные поля
- 5 Задача классификации последовательности
  - Условные случайные поля
  - Рекуррентные нейронные сети
  - Заполнение слотов

# Условные случайные поля [Lafferty, 2001]

## Условные случайные поля [Conditional random fields]

$$\hat{Y} = \arg \max_Y P(Y|X) = \phi(y_i, y_{i-1})\phi(y_i, x_i)$$



Источник: [http://davidsbatista.net/blog/2017/11/13/Conditional\\_Random\\_Fields/](http://davidsbatista.net/blog/2017/11/13/Conditional_Random_Fields/)



# Условные случайные поля

Вероятность последовательности меток классов для входной последовательности определяется по признакам, которые называются потенциальными функциями. Эти признаки помогают связать класс текущего наблюдения  $x_i$  с классами других наблюдений. Для формализации признаков чаще всего используются индикаторные функции. Таким образом, задача обучения сводится к определению весов индикаторных функций.

$$t(y_{i-1}, y_i, x, i) = \begin{cases} 1, & \text{if } x_i = \text{"June"} \text{ and } y_{i-1} = IN \text{ and } y_i = NNP \\ 0, & \text{otherwise} \end{cases}$$

$$s(y_{i-1}, x, i) = \begin{cases} 1, & \text{if } x_i = \text{"to"} \text{ and } y_i \\ 0, & \text{otherwise} \end{cases}$$

# Условные случайные поля

Для того, чтобы найти вероятность последовательности классов для входной последовательности:

- 1 извлекаем признаки
- 2 находим их веса и линейную комбинацию их признаков с найденными весами
- 3 используем softmax для определения искомых вероятностей.

Обозначим все признаки:  $f(y_{i-1}, y_i, x, i)$ . Признаки для последовательностей:  $F(y, x) = \sum_{i=1}^n f(y_{i-1}, y_i, x, i)$ . Обозначим веса признаков через  $\lambda$ . Искомая вероятность:

$$p(y|x) = \frac{e^{\sum_{i=1}^k \lambda_i F_i(y, x)}}{\sum_{y' \in C^n} e^{\sum_{i=1}^k \lambda_i F_i(y', x)}}$$

## Пример. NER

Город/О Пушкин/(Per, City) является/О ...

Признаки:

$$f_1(y_{i-1}, y_i, x, i) = \begin{cases} 1, & \text{if } x_i = \text{Пушкин and } y_{i-1} = O \text{ and } y_i = \text{Per} \\ 0, & \text{otherwise} \end{cases}$$

$$f_2(y_{i-1}, y_i, x, i) = \begin{cases} 1, & \text{if } x_i = \text{Пушкин and } y_{i-1} = O \text{ and } y_i = \text{City} \\ 0, & \text{otherwise} \end{cases}$$

$$f_3(y_{i-1}, y_i, x, i) = \begin{cases} 1, & \text{if } x_i = \text{Пушкин and } x_{i+1}[:3] = \text{тс} \text{ and } y_i = \text{Per} \\ 0, & \text{otherwise} \end{cases}$$

$$f_4(y_{i-1}, y_i, x, i) = \begin{cases} 1, & \text{if } x_i = \text{Пушкин and } x_{i+1}[:3] = \text{тс} \text{ and } y_i = \text{City} \\ 0, & \text{otherwise} \end{cases}$$

$$f_5(y_{i-1}, y_i, x, i) = \begin{cases} 1, & \text{if } x_i = \text{Пушкин and } x_{i-1} = \text{город and } y_i = \text{Per} \\ 0, & \text{otherwise} \end{cases}$$

$$f_6(y_{i-1}, y_i, x, i) = \begin{cases} 1, & \text{if } x_i = \text{Пушкин and } x_{i-1} = \text{город and } y_i = \text{City} \\ 0, & \text{otherwise} \end{cases}$$

Веса:  $\lambda_1 = 5, \lambda_2 = 2, \lambda_3 = 10, \lambda_4 = 7, \lambda_5 = 7, \lambda_6 = 20$ .

## Пример. NER

Сравним вероятности  $P(O \text{ Per } O \mid \text{Город Пушкин является})$  и  $P(O \text{ City } O \mid \text{Город Пушкин является})$ .

①  $O \text{ Per } O : 5 + 10 + 7 = 12$

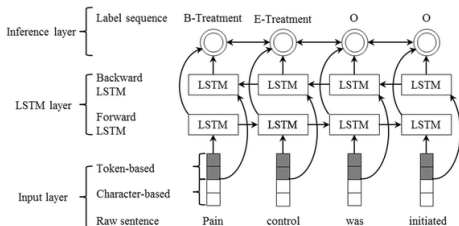
②  $O \text{ City } O : 2 + 7 + 20 = 27$

$$P(O \text{ Per } O \mid \text{Город Пушкин является}) = \frac{e^{12}}{e^{12} + e^{27}}$$

$$P(O \text{ City } O \mid \text{Город Пушкин является}) = \frac{e^{27}}{e^{12} + e^{27}}$$

- 1 Извлечение информации
- 2 Извлечение информации по словарям, регулярными выражениями и по шаблонам
- 3 Классификация чанков [chunk]
- 4 Условные случайные поля
- 5 Задача классификации последовательности
  - Условные случайные поля
  - Рекуррентные нейронные сети
  - Заполнение слотов

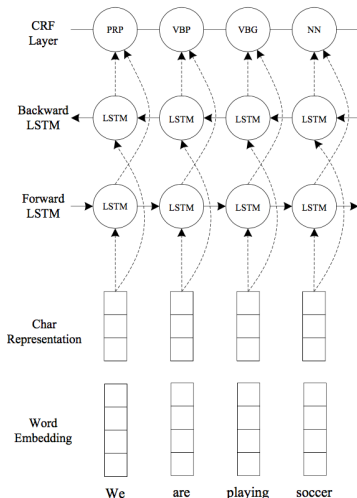
# biLSTM для классификации последовательности



Каждый выход нейронной сети решает свою задачу классификации: какую метку приписать входному слову?

Источник: <https://bmcmedinformdecismak.biomedcentral.com/articles/10.1186/s12911-017-0468-7>

# (CNN-)biLSTM-CRF [Ma and Hovy]



Слой с CRF используется для перевзвешивания выхода нейронной сети:

$$\begin{aligned} \text{score}_{lstm-crf}(x, y) &= \\ &= \sum_{i=0}^n W_{y_{i-1}, y_i} \cdot \text{LSTM}(x_i) + b(y_{i-1}, y_i) \end{aligned}$$

- 1 Извлечение информации
- 2 Извлечение информации по словарям, регулярными выражениями и по шаблонам
- 3 Классификация чанков [chunk]
- 4 Условные случайные поля
- 5 Задача классификации последовательности
  - Условные случайные поля
  - Рекуррентные нейронные сети
  - Заполнение слотов



# Меры качества классификации последовательностей

## 1 token-based

$tp$  – число истинно-положительных токенов,  $fp$  – число ложно-положительных токенов,  $fn$  – число ложно-отрицательных токенов

## 2 chunk-based

чанк – именованная сущность (синтаксическая группа, и др.) целиком

$tp$  – число истинно-положительных чанков,  $fp$  – число ложно-положительных чанков,  $fn$  – число ложно-отрицательных чанков