

# Задача поиска символов в текстах

Севериков Павел Андреевич

Московский физико-технический институт

*Курс:* Численные методы обучения по прецедентам  
(практика, В. В. Стрижов)/Группа 674, весна 2019

## Проблема

Современные модели для обработки текстов воспринимают высказывания буквально, и различные средства художественной выразительности, в частности, символы, метафоры, аллегории и др. не интерпретируются ими верным образом.

## Цель работы

Получить оптимальную модель для определения неоднозначности в высказываниях.

## Золотые руки

- "Мастер с золотыми руками"— умение хорошо что-либо делать
- "У статуи золотые руки"— материал, из которого сделана статуя

-  Marek Rei, Gamal K.O. Crichton, Sampo Pyysalo N. Attending to Characters in Neural Sequence Labeling Models // Proceedings of COLING 2016, the 26th International Conference on Computational Linguistics: Technical Papers , 2016, C16-1030, Pp. 309–318.
-  Adnan Akhundov, Dietrich Trautmann, Georg Groh N. Sequence Labeling: A Practical Approach // CoRR , vol. abs/1808.03926, 2018.
-  Gao, Ge and Choi, Eunsol and Choi, Yejin and Zettlemoyer, Luke N. Neural Metaphor Detection in Context // Proceedings of the 2018 Conference on Empirical Methods in Natural Language Processing, 2018, Pp. 607–613.

## Sequence labeling

Дано предложение  $\mathbf{X}$ , разделенное на слова:  $\{x_1, x_2, \dots, x_n\}$ .  
Требуется построить последовательность двоичных меток (labels)  $\{l_1, l_2, \dots, l_n\}$ , которые идентифицируют наличие неоднозначности/символа в каждом слове  $x_i$

## Классификация

Требуется для целевой переменной  $i$  предсказать отношение  $x_i$  к классу символ или не символ, соответственно 1 и 0.

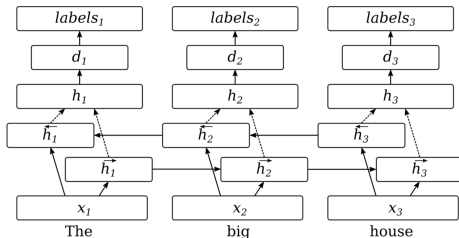


Рис.: Схема базовой модели BiLSTM

## Представления в LSTM-сети

$$\vec{h}_t = \text{LSTM}(x_t, \vec{h}_{t-1}) \quad \overleftarrow{h}_t = \text{LSTM}(x_t, \overleftarrow{h}_{t+1}) \quad h_t = [\vec{h}_t; \overleftarrow{h}_t]$$

Скрытый слой нелинейности:  $d_t = \tanh(W_d h_t)$ , где  $W_d$  – весовая матрица между слоями.

# Итоговая задача оптимизации

Нормированное распределение вероятностей по всем возможным меткам для каждого слова (softmax):

$$\mathbb{P}(y_t = k | d_t) = \frac{e^{W_k d_t}}{\sum_{\tilde{k} \in K} e^{W_{\tilde{k}} d_t}},$$

где  $\mathbb{P}(y_t = k | d_t)$  – вероятность того, что метка  $t$ -ого слова  $y_t$  будет  $k$  ( $K$  – множество всевозможных меток),  $W_k$  –  $k$ -ая строка весовой матрицы  $W$ .

Для оптимизации модели используется минимизация функции

$$\mathcal{L} = - \sum_{t=1}^T \log(\mathbb{P}(y_t | d_t))$$

Т.е. решается данная задача:

$$W^* = \operatorname{argmin}_W (\mathcal{L}(W))$$

## Тестируемая модель

BiLSTM нейронная сеть с softmax Возможные улучшения:

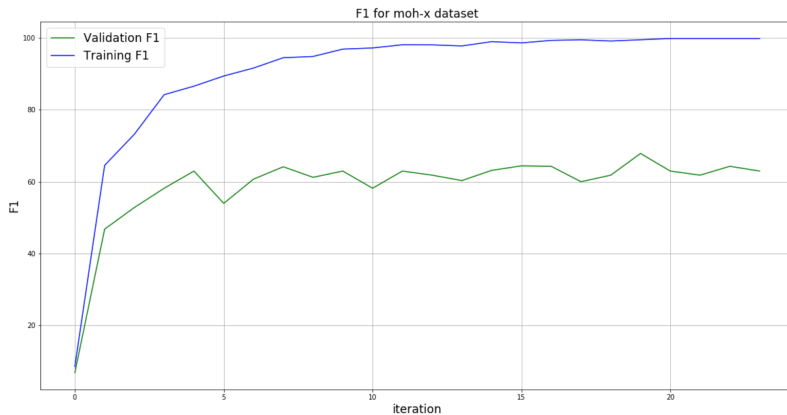
- 1 BiLSTM нейронная сеть с CRF (Conditional random field)
- 2 BiLSTM нейронная сеть с Attention

## Данные

- 1 MOH датасет с метафорами (eng)
- 2 VU Amsterdam Metaphor Corpus (eng)
- 3 Разметка для текста "Мастер и Маргарита"
- 4 Размеченные данные из института русского языка



# Результаты эксперимента: F1 мера для МОН датасета



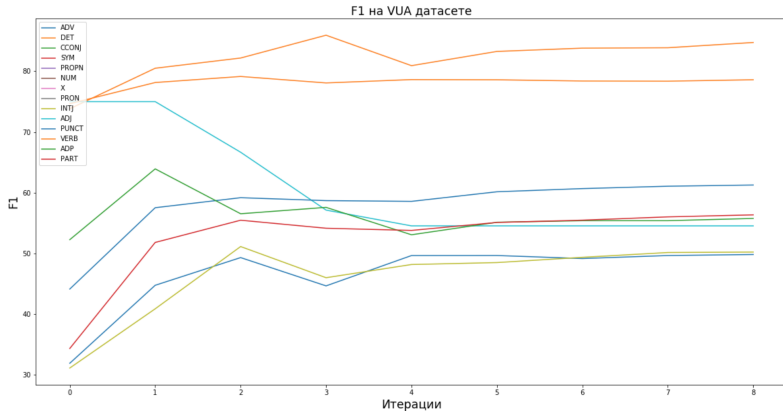
После 10 эпох обучения на МОН датасете:

- Precision = **64.14**
- Recall = **67.86**
- F1 = **65.89**
- Accuracy = **69.27**

После 10 эпох обучения на датасете "Мастер и Маргарита":

- Precision = **86.96**
- Recall = **93.02**
- F1 = **89.89**
- Accuracy = **88.61**

# Результаты эксперимента: F1 для VUA датасета:



- Алгоритм sequence labeling хорошо подходит для поиска символов в тексте
- Качество заметно улучшится при увеличении выборки
- Предложенные модели могут быть применены для определения не только каких-то конкретных неоднозначностей в тексте, а в целом для всех видов символов
- Для русскоязычных текстов данная задача никак до этого не решалась