Основы обработки текстов

Лекция #6: Базовые задачи обработки текстов

Задача NERC

- На входе: текст, разбитый на предложения и токены
- На выходе: множество сущностей (подстрока, текст)

Microsoft — один из крупнейших производителей ПО в мире

Задача NERC

- На входе: текст, разбитый на предложения и токены
- На выходе: множество сущностей (подстрока, текст)

Александр Пушкин родился в Москве, столице России **Александр Пушкин** родился в <mark>Москве</mark>, столице России

[PER]

[LOC]

[LOC]

Microsoft — один из крупнейших производителей ПО в мире



Microsoft — один из крупнейших производителей ПО в мире

Сегментация текста

- На входе: текст (последовательность символов)
- На выходе: текст, разбитый на предложения и токены

Александр Пушкин родился в Москве, столице России



Александр Пушкин родился в Москве, столице России

Microsoft — один из крупнейших производителей ПО в мире



Microsoft — один из крупнейших производителей ПО в мире

Сегментация текста

- На входе: текст (последовательность символов)
- На выходе: текст, разбитый на предложения и токены
 - Предложение это единица языка, которая представляет собой грамматически организованное соединение слов, обладающее смысловой законченностью.

• Токен – минимальная лингвистическая единица текста (речи), обладающая смыслом.

{слово, пунктуация, число}

- На входе: текст (последовательность символов)
- На выходе: текст, разбитый на предложения и токены

Александр Пушкин родился в Москве, столице России.

• Решение (RegExp, регулярные выражения)

- На входе: текст (последовательность символов)
- На выходе: текст, разбитый на предложения и токены

Александр Пушкин родился в Москве, столице России.

• Решение (RegExp, регулярные выражения)

Разделить входную строку по пробельным символам

Александр Пушкин родился в Москве, столице России.

- На входе: текст (последовательность символов)
- На выходе: текст, разбитый на предложения и токены

Александр Пушкин родился в Москве, столице России.

• Решение (RegExp, регулярные выражения)

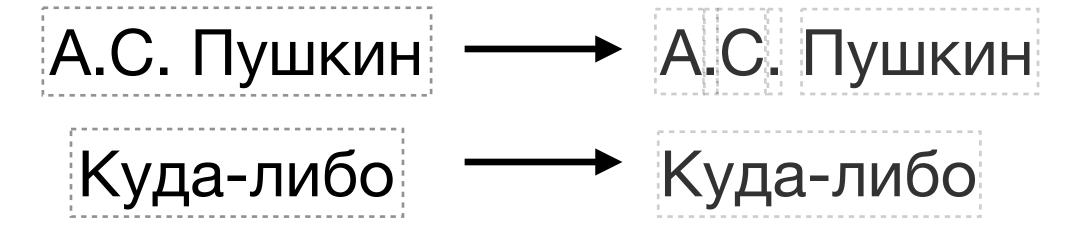
Разделить входную строку по пробельным символам

Александр Пушкин родился в Москве, столице России.

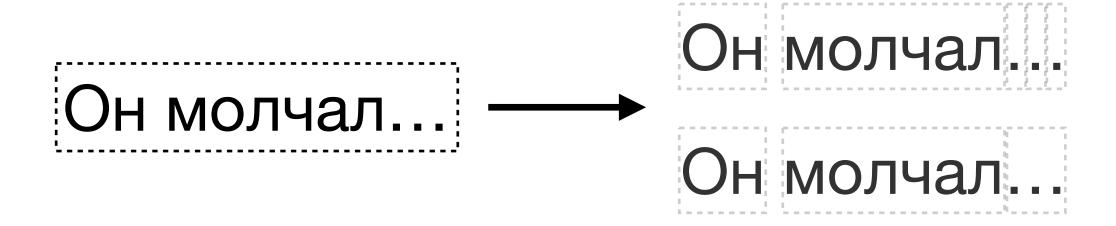
Отделить префиксные и постфиксные знаки пунктуации

Александр Пушкин родился в Москве, столице России.

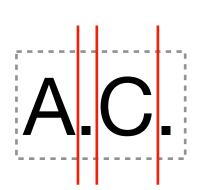
• Пунктуация бывает внутри слова

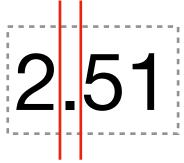


• Может быть несколько символов пунктуации подряд



- Будем считать токенами подстроки, содержащие только буквы и цифры (alphanumeric)
- Остальные подстроки попробуем разделить по не alphanumeric символам, используя бинарный классификатор







Признаки:

- Длины токеновкандидатов
- Символы (n-gram)
- Словарь
- . . .

Классификаторы:

- SVM
- Лог. регрессия
- Нейронные сети
- •

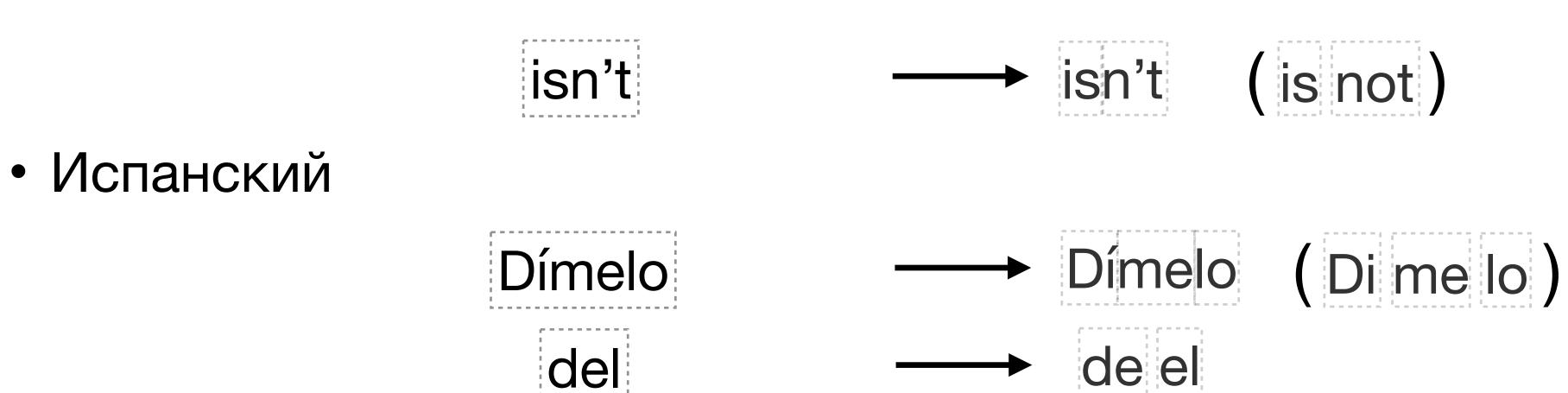
- Некоторые токены должны состоять из нескольких слов
 - Русский

Как только началось ---- Как только началось

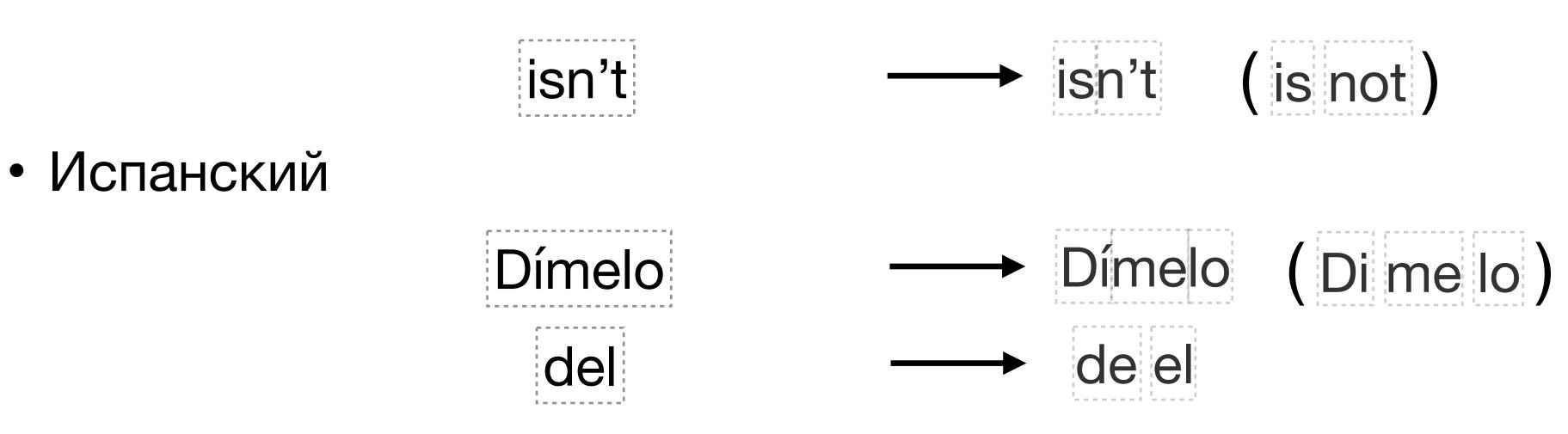
• Испанский

¿Cómo te llamas? → ¿Cómo te llamas?

- Некоторые слова склеиваются
 - Английский



- Некоторые слова склеиваются
 - Английский



• Немецкий

Rechtsschutzversicherungsgesellschaften

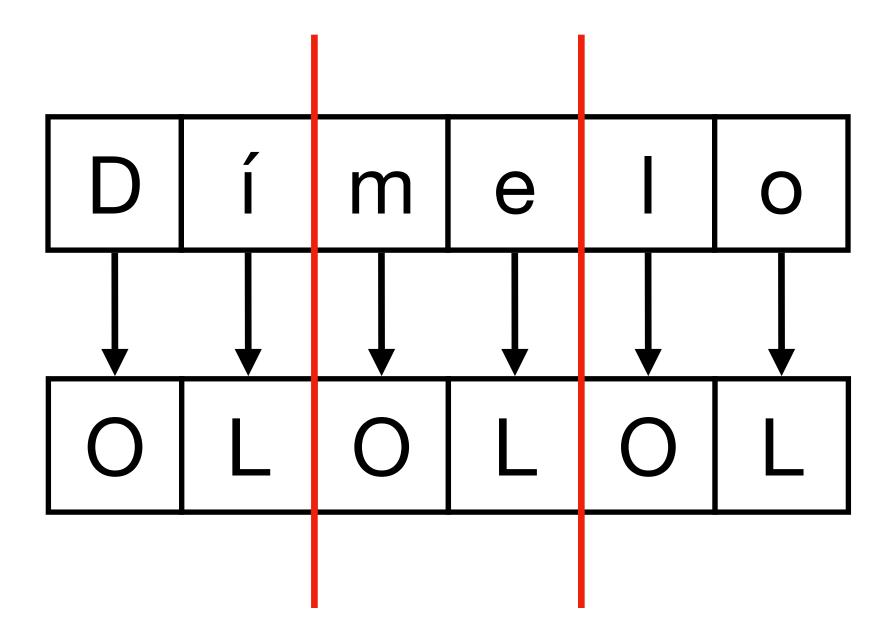
• Задача разметки последовательности. Каждый символ необходимо

отнести к одному из двух классов:

• Последний символ в токене (L)

• Не последний символ в токене (О)

- Методы:
 - RNN (Рекурентные нейронные сети)
 - BiRNN
 - •



• После разделения на символы замена токенов по словарю

Dímelo → (Di me lo)

Определение границ предложений

На входе:

- текст (последовательность токенов)
- текст (последовательность символов)

На выходе: текст, разбитый на предложения

Идея:

• В естественных языках предложения отделяются друг от друга знаками препинания (.!?) (。в китайском) (: в армянском)

Решение:

• Классифицировать каждый токен (.!?) является ли он концом предложения

Сегментация текста

На входе: текст (последовательность символов)

На выходе: текст, разбитый на предложения и токены

Идея

Одновременно находить границы токенов и предложений

Решение

- Разметка последовательности. Каждый символ классифицируется на один из **трех** классов:
 - Последний символ предложения
 - Последний символ токена
 - Обычный символ

Оценка качества

Gold

Александр Пушкин родился в Москве, столице России.

Predicted

Александр Пушкин родился в Москве, столице России.

Оценка качества

Gold

Александр Пушкин родился в Москве, столице России.

Predicted

Александр Пушкин родился в Москве, столице России.

• Метрики

$$P = \frac{|\text{correct}|}{|\text{predicted}|}$$

$$R = \frac{|\text{correct}|}{|\text{gold}|}$$

$$R = \frac{5}{9}$$

$$F_1 = \frac{2 \cdot P \cdot R}{P + R}$$

- На входе: текст (последовательность символов)
- На выходе: метка языка

Александр Сергеевич Пушкин родился в Москве 26 мая 1799 года.

Александр Сергей улы Пушкин 1799 елның 26 маенда Мәскәү шәһәрендә туа.

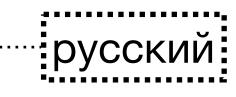
Аляксандр Сяргеевіч Пушкін нарадзіўся ў Маскве 26 мая 1799 года.

Alexander Sergej Pushkin fæddist í Moskvu 26 Maí 1799.

ალექსანდრე სერგეევიჩ პუშკინი დაიბადა მოსკოვში 1799 წლის 26 მაისს.

- На входе: текст (последовательность символов)
- На выходе: метка языка

Александр Сергеевич Пушкин родился в Москве 26 мая 1799 года.



Александр Сергей улы Пушкин 1799 елның 26 маенда Мәскәү шәһәрендә туа.

Аляксандр Сяргеевіч Пушкін нарадзіўся ў Маскве 26 мая 1799 года.

Alexander Sergej Pushkin fæddist í Moskvu 26 Maí 1799.

ალექსანდრე სერგეევიჩ პუშკინი დაიბადა მოსკოვში 1799 წლის 26 მაისს.

- На входе: текст (последовательность символов)
- На выходе: метка языка

Александр Сергеевич Пушкин родился в Москве 26 мая 1799 года.

Александр Сергей улы Пушкин 1799 елның 26 маенда Мәскәү шәһәрендә туа.

Аляксандр Сяргеевіч Пушкін нарадзіўся ў Маскве 26 мая 1799 года.

Alexander Sergej Pushkin fæddist í Moskvu 26 Maí 1799.

ალექსანდრე სერგეევიჩ პუშკინი დაიბადა მოსკოვში 1799 წლის 26 მაისს.

русский

татарский

- На входе: текст (последовательность символов)
- На выходе: метка языка

Александр Сергеевич Пушкин родился в Москве 26 мая 1799 года.

Александр Сергей улы Пушкин 1799 елның 26 маенда Мәскәү шәһәрендә туа.

Аляксандр Сяргеевіч Пушкін нарадзіўся ў Маскве 26 мая 1799 года.

Alexander Sergej Pushkin fæddist í Moskvu 26 Maí 1799.

ალექსანდრე სერგეევიჩ პუშკინი დაიბადა მოსკოვში 1799 წლის 26 მაისს.

русский

татарский

белорусский

- На входе: текст (последовательность символов)
- На выходе: метка языка

Александр Сергеевич Пушкин родился в Москве 26 мая 1799 года.

Александр Сергей улы Пушкин 1799 елның 26 маенда Мәскәү шәһәрендә туа.

Аляксандр Сяргеевіч Пушкін нарадзіўся ў Маскве 26 мая 1799 года.

Alexander Sergej Pushkin fæddist í Moskvu 26 Maí 1799.

ალექსანდრე სერგეევიჩ პუშკინი დაიბადა მოსკოვში 1799 წლის 26 მაისს.

русский

татарский

белорусский

исландский

მოსკოვში 1799 წლის 26 მაისს.

- На входе: текст (последовательность символов)
- На выходе: метка языка

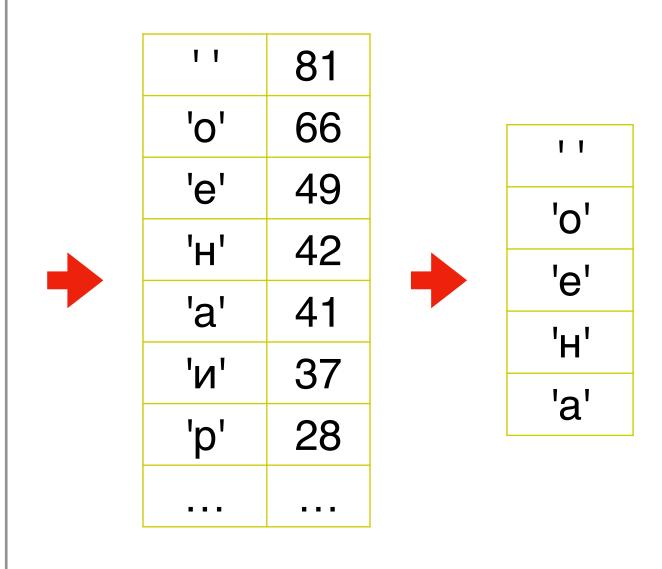
Александр Сергеевич Пушкин родился в Москве русский 26 мая 1799 года. Александр Сергей улы Пушкин 1799 елның 26 маенда Мәскәү шәһәрендә туа. Аляксандр Сяргеевіч Пушкін нарадзіўся ў белорусский Маскве 26 мая 1799 года. Alexander Sergej Pushkin fæddist í Moskvu 26 Maí исландский 1799. ალექსანდრე სერგეევიჩ პუშკინი დაიბადა [грузинский]

Профили языка*

Для каждого языка извлекается его профиль

• Профиль – список (top-K) символьных n-gram отсортированный по невозрастанию частоты

Происхождение Александра Сергеевича Пушкина идёт от разветвлённого нетитулованного дворянского рода Пушкиных, восходившего по генеалогической легенде к «мужу честну» Ратше. Пушкин неоднократно писал о своей родословной в стихах и прозе; он видел в своих предках образец истинной «аристократии», древнего рода, честно служившего отечеству, но не снискавшего благосклонности правителей и «гонимого». Не раз он обращался (в том числе в художественной форме) и к образу своего прадеда по матери — африканца Абрама Петровича Ганнибала, ставшего слугой и воспитанником Петра I, а потом военным инженером и генералом.



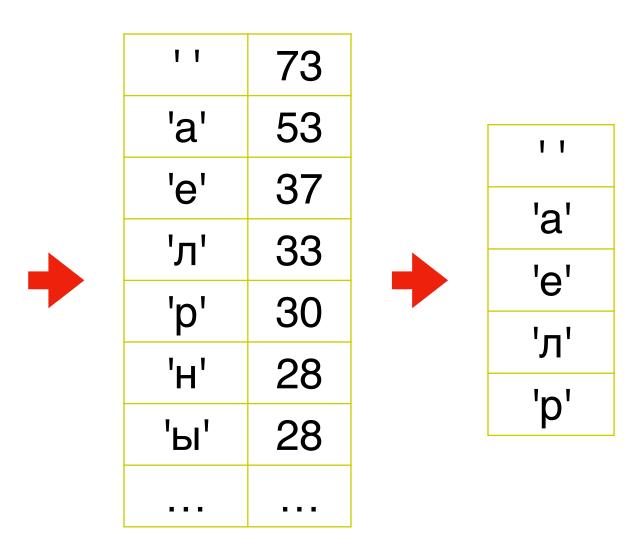
^{*} W. B. Cavnar, J. M. Trenkle. "N-gram-based text categorization." (1994)

Профили языка

Для каждого языка извлекается его профиль

• Профиль – список (top-K) символьных n-gram отсортированный по невозрастанию частоты

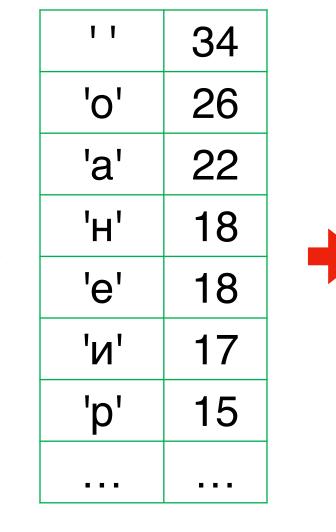
Үсмер чагының 6 елын Пушкин 1811 елның 19 октябрендә дәрәҗәле кешеләр өчен ачылган лицейда үткәрә. Лицейда укыгында усмернең таланты ачыла һәм башкалар тарафыннан югары бәяләнә. Лицейда үткәрелгән еллар, дуслары шагыйры күңелендә мәңгегә саклана һәм иҗатында чагылыш таба. Александр Сергеевичның яшьлеге шушында үтә, биредә шагыйры Пушкин туа: 130лап шигыры иҗат ителә, «Руслан һәм Людмила» поэмасы языла башлый, танылган журналларда беренче язмалары басыла. Лицей хәзерге урта яки югары уку йортлары дәрәҗәсендә белем бирергә тиеш була.



Профили языка

Для классифицируемого текста извлекается профиль

Весной 1820 года Пушкина вызвали к военному генералгубернатору Петербурга графу М. А. Милорадовичу для объяснения по поводу содержания его стихотворений (в том числе эпиграмм на Аракчеева, архимандрита Фотия и самого Александра I), несовместимых со статусом государственного чиновника.



1 1

'0'

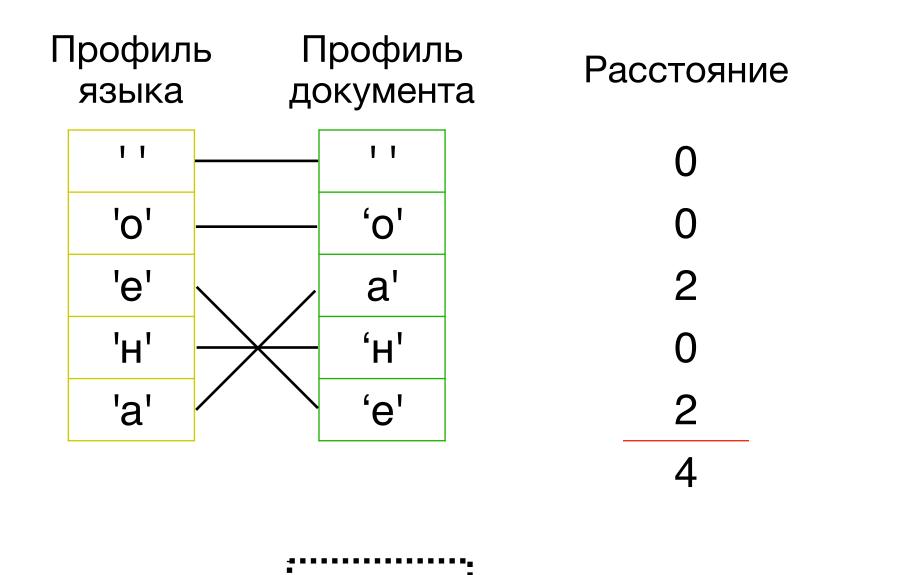
'a'

Ή'

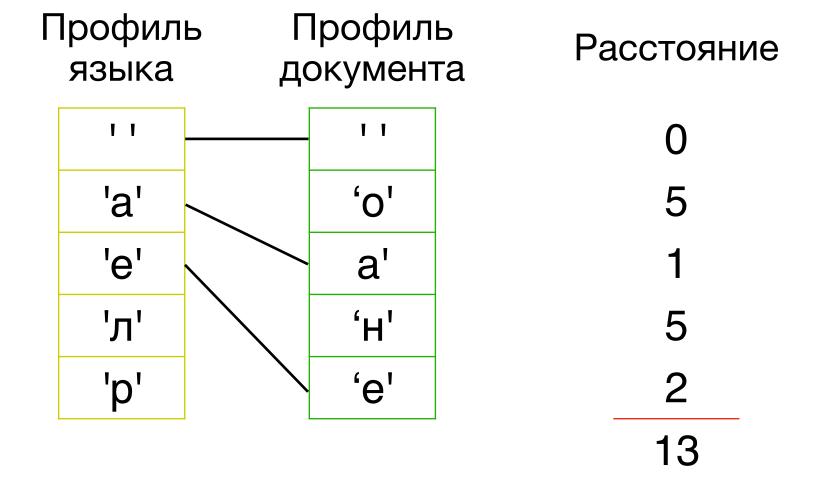
'e'

Профили языка

Профиль классифицируемого текста сравнивается с профилями языков (out-of-place), выбирается ближайший язык



русский



татарский

• Теорема Байеса

$$p(c \mid d) = \frac{p(d \mid c) \cdot p(c)}{p(d)}$$

- $p(c \mid d)$ вероятность $d \in c$
- $p(d \mid c)$ вероятность встретить d в c
- p(c) вероятность встретить документ класса c
- p(d) вероятность встретить документ d
- Классификатор

$$\hat{c} = \arg\max_{c \in C} p(c \mid d) = \arg\max_{c \in C} \frac{p(d \mid c) \cdot p(c)}{p(d)} = \arg\max_{c \in C} p(d \mid c) \cdot p(c)$$

$$\hat{c} = \arg \max_{c \in C} p(d \mid c) \cdot p(c)$$

• Представим d как $[f_1, f_2, ..., f_n]$ (набор признаков)

$$p(d | c) = p(f_1, f_2, ..., f_n | c) = p(f_1 | c) \cdot p(f_2, ..., f_n | c, f_1) =$$

$$= p(f_1 | c) \cdot p(f_2 | c, f_1) \cdot ... \cdot p(f_n | c, f_1, ..., f_{n-1})$$

$$\hat{c} = \arg \max_{c \in C} p(d \mid c) \cdot p(c)$$

• Представим d как $[f_1, f_2, ..., f_n]$ (набор признаков)

$$p(d | c) = p(f_1, f_2, ..., f_n | c) = p(f_1 | c) \cdot p(f_2, ..., f_n | c, f_1) =$$

$$= p(f_1 | c) \cdot p(f_2 | c, f_1) \cdot ... \cdot p(f_n | c, f_1, ..., f_{n-1})$$

• Предположим, что f_i независимы при условии c

$$p(f_i | c, f_j) = p(f_i | c), i \neq j$$

$$\hat{c} = \arg \max_{c \in C} p(d \mid c) \cdot p(c)$$

• Представим d как $[f_1, f_2, ..., f_n]$ (набор признаков)

$$p(d|c) = p(f_1|c) \cdot p(f_2|c) \cdot \dots \cdot p(f_n|c) = \prod_{i=1}^{n} p(f_i|c)$$

• Предположим, что f_i независимы при условии c

$$p(f_i | c, f_j) = p(f_i | c), i \neq j$$

Naïve Bayes Classifier

Теория

$$\hat{c} = \arg \max_{c \in C} p(d \mid c) \cdot p(c)$$

• Представим d как $[f_1, f_2, ..., f_n]$ (набор признаков)

$$\hat{c} = \arg \max_{c \in C} \prod_{i=1}^{n} p(f_i | c) \cdot p(c)$$

$$\hat{c} = \arg \max_{c \in C} \left[\sum_{i=1}^{n} \log p(f_i | c) + \log p(c) \right]$$

Определение языка текста

Naïve Bayes Classifier

$$\hat{c} = \arg\max_{c \in C} \left[\sum_{i=1}^{n} \log p(f_i | c) + \log p(c) \right]$$

- Признаки символьные n-gram'ы
- Оценка параметров:

$$p(c) = \frac{|D_c|}{|D|};$$
 или $p(c) = \frac{1}{|C|};$
$$p(f_i|c) = \frac{\left|\{f_i \in D_c\}\right|}{\sum_j \left|\{f_j \in D_c\}\right|}$$

Naïve Bayes Classifier

Сглаживание

• Проблема: некоторые n-gram'ы могут не встретиться в обучающем корпусе

$$p(f_i | c) = \frac{\left| \{ f_i \in D_c \} \right|}{\sum_j \left| \{ f_j \in D_c \} \right|} = \frac{0}{\sum_j \left| \{ f_j \in D_c \} \right|} = 0$$

• Решение: аддитивное сглаживание

$$p(f_i | c) = \frac{\left| \{ f_i \in D_c \} \right| + \varepsilon}{\sum_j \left(\left| \{ f_j \in D_c \} \right| + \varepsilon \right)} \neq 0$$

Naïve Bayes Classifier

Обучение

Происхождение Александра Сергеевича Пушкина идёт от разветвлённого нетитулованного дворянского рода Пушкиных, восходившего по генеалогической легенде к «мужу честну»

:русский

Ратш Стиха Дережеле кешелер ечен ачылган лицейда үткере. Лицейда «арис укыгында усмернең таланты ачыла һәм башкалар тарафыннан но не не ра күңеленде мәңгеге саклана һәм иҗатында чагылыш таба. Александр Сергеевичның яшьлеге шушында үтә, биреде Петок жәрелген алан һәм ижатында чагылыш таба. Петр шагыйрь Пушкин туа: 130лап шигырь иҗат ителә, «Руслан һәм а пот Людмила» поэмасы языла башлый, танылган журналларда беренче язмалары басыла. Лицей хәзерге урта яки югары уку йортлары дәрәҗәсендә белем бирергә тиеш була.

$$p$$
(русский) = $\frac{1}{2}$

$$p$$
(татарский) = $\frac{1}{2}$

Обучение



Происхождение Александра Сергеевича Пушкина идёт от разветвлённого нетитулованного дворянского рода Пушкиных, восходившего по генеалогической легенде к «мужу честну» Ратше. Пушкин неоднократно писал о своей родословной в стихах и прозе; он видел в своих предках образец истинной «аристократии», древнего рода, честно служившего отечеству, но не снискавшего благосклонности правителей и «гонимого». Не раз он обращался (в том числе в художественной форме) и к образу своего прадеда по матери — африканца Абрама Петровича Ганнибала, ставшего слугой и воспитанником Петра I, а потом военным инженером и генералом.

$$p(\Pi \mid \text{русский}) = \frac{6+1}{610+46} = 0.01067 \qquad p(\text{и} \mid \text{русский}) = \frac{37+1}{610+46} = 0.05793$$

$$p(\text{р} \mid \text{русский}) = \frac{28+1}{610+46} = 0.04421 \qquad p(\text{с} \mid \text{русский}) = \frac{27+1}{610+46} = 0.04268 \qquad \dots$$

$$p(\text{о} \mid \text{русский}) = \frac{66+1}{610+46} = 0.10213 \qquad p(\text{ж} \mid \text{русский}) = \frac{5+1}{610+46} = 0.00914$$

Обучение



Үсмер чагының 6 елын Пушкин 1811 елның 19 октябрендә дәрәҗәле кешеләр өчен ачылган лицейда үткәрә. Лицейда укыгында усмернең таланты ачыла һәм башкалар тарафыннан югары бәяләнә. Лицейда үткәрелгән еллар, дуслары шагыйрь күңелендә мәңгегә саклана һәм иҗатында чагылыш таба. Александр Сергеевичның яшьлеге шушында үтә, биредә шагыйрь Пушкин туа: 130лап шигырь иҗат ителә, «Руслан һәм Людмила» поэмасы языла башлый, танылган журналларда беренче язмалары басыла. Лицей хәзерге урта яки югары уку йортлары дәрәҗәсендә белем бирергә тиеш була.

$$p(\Pi \mid \text{русский}) = \frac{2+1}{538+54} = 0.00507 \qquad p(\text{и} \mid \text{русский}) = \frac{16+1}{538+54} = 0.02872$$

$$p(\text{р} \mid \text{русский}) = \frac{30+1}{538+54} = 0.05236 \qquad p(\text{с} \mid \text{русский}) = \frac{9+1}{538+54} = 0.01689 \qquad \dots$$

$$p(\text{о} \mid \text{русский}) = \frac{3+1}{538+54} = 0.00675 \qquad p(\text{ж} \mid \text{русский}) = \frac{1+1}{538+54} = 0.00337$$

Предсказание

Весной 1820 года Пушкина вызвали к военному генерал-губернатору Петербурга графу М. А. Милорадовичу для объяснения по поводу содержания его стихотворений (в том числе эпиграмм на Аракчеева, архимандрита Фотия и самого Александра I), несовместимых со статусом государственного чиновника.

$$\hat{c} = \arg\max_{c \in C} \left[\sum_{i=1}^{n} \log p(f_i | c) + \log p(c) \right]$$

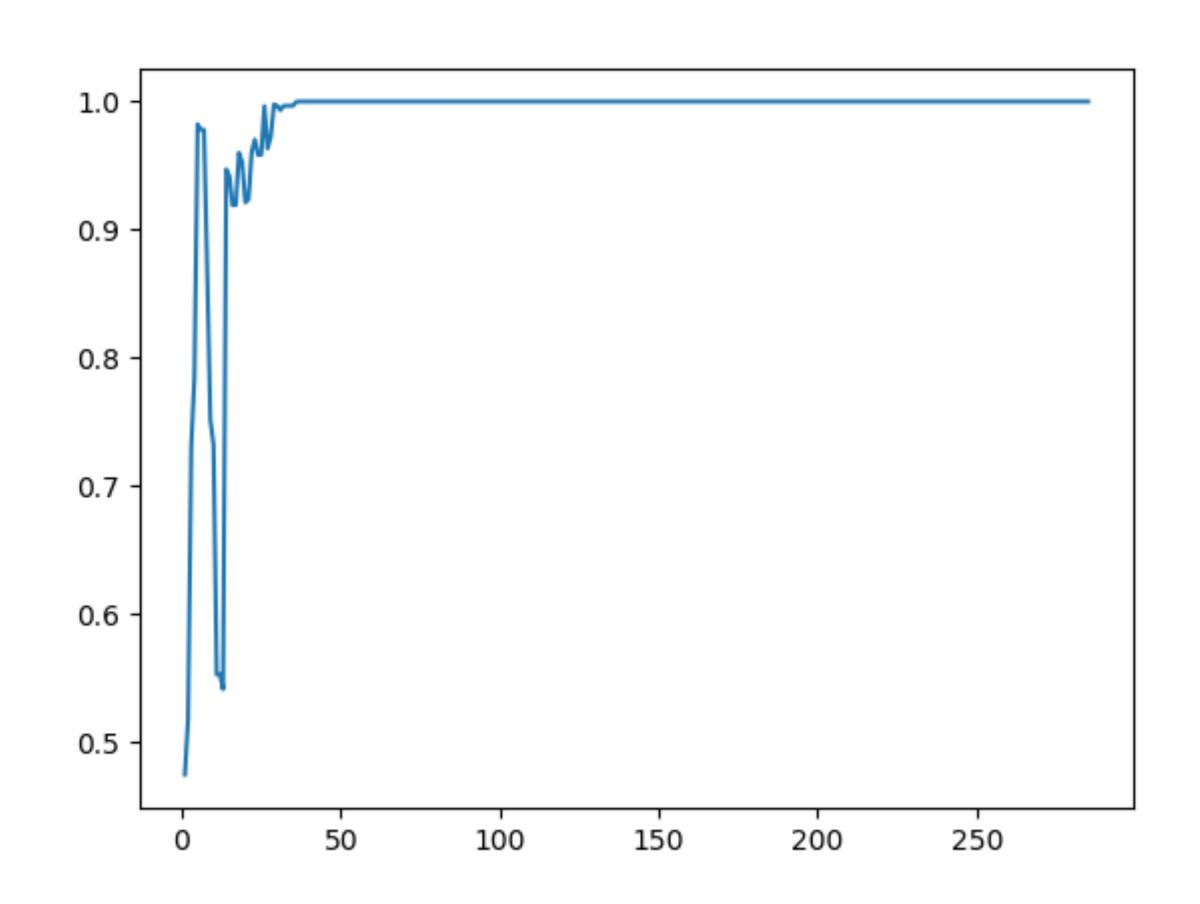
```
\log p(c) \quad \log p(\mathsf{B} \, | \, c) \quad \log p(\mathsf{e} \, | \, c) \quad \log p(\mathsf{c} \, | \, c) \hat{c}_{\mathsf{русский}} = -0.69315 - 6.48768 - 2.57566 - 3.15547 - \ldots - 5.10139 \hat{c}_{\mathsf{русский}} = -953.05578
```

Предсказание

• Назад к вероятностям

$$p(c_i|d) = \frac{\exp \hat{c}_i}{\sum_{j} \exp \hat{c}_j} = \frac{1}{1 + \sum_{j \neq i} \exp \hat{c}_j - \hat{c}_i}$$

- Наш пример
 - $\hat{c}_{\text{русский}} = -953.05578;$ $\hat{c}_{\text{татарский}} = -1066.1$
 - $p(\text{русский} \mid d) = \frac{1}{1 + e^{-113.12848}} \approx 1$
- Для первых 30 символов
 - $\hat{c}_{\text{русский}} = -113.25854$; $\hat{c}_{\text{татарский}} = -118.98256$
 - . p(русский $|d) = \frac{1}{1 + e^{-5.72402}} \approx 0.99674$



Весной 1820 года Пушкина вызва

- Задача многоклассовой классификации (multiclass classification)
 - Объекты тексты
 - Классы метки языков

- Признаки классификации
 - N-gram по символам
 - Слова (из словаря)

- Классификаторы
 - Multiclass SVM
 - Naïve Bayes
 - Нейронные сети
 - Perceptron
 - Convolutional NN)

•

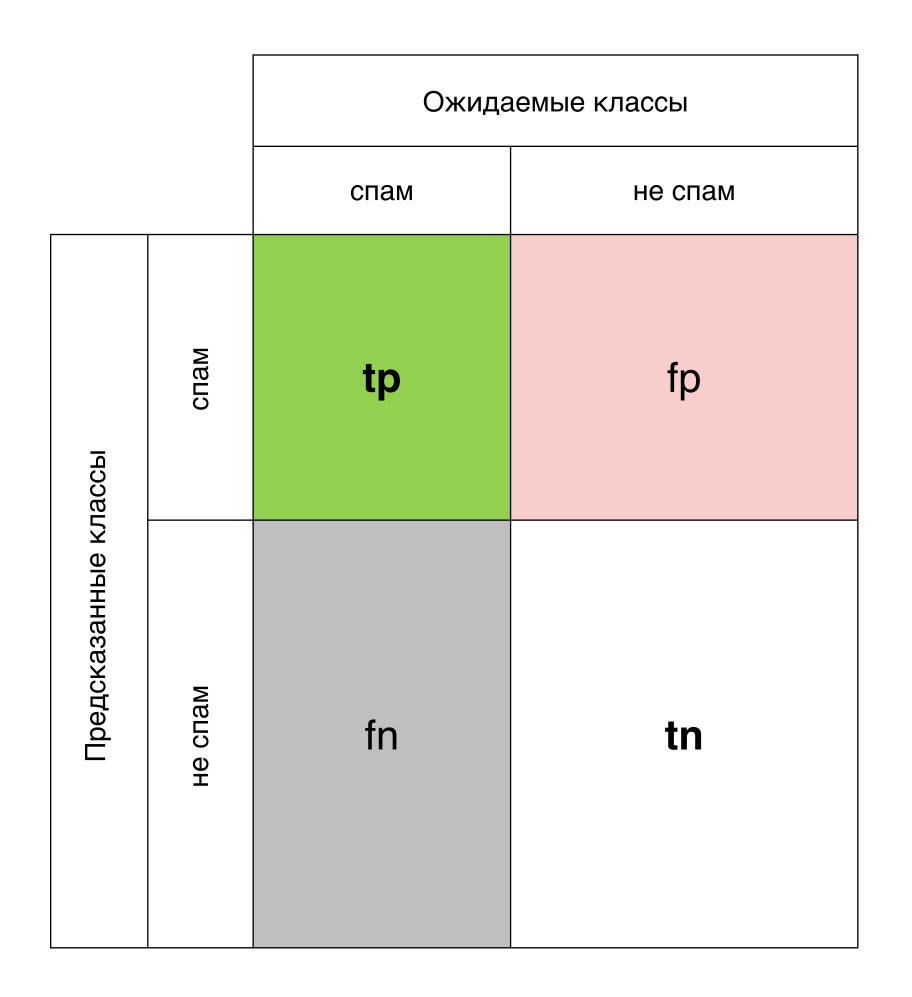
Binary classification

Accuracy =
$$\frac{tp + tn}{tp + tn + fp + fn}$$

$$Precision = \frac{tp}{tp + fp}$$

$$Recall = \frac{tp}{tp + fn}$$

$$F_1 = \frac{2 \cdot Precision \cdot Recall}{Precision + Recall}$$



Confusion table

Оценка качестваmulticlass classification

| | | Ожидаемые классы | | | | |
|----------------------|------------|------------------|-----------|------------|--|--|
| | | русский | татарский | английский | | |
| Предсказанные классы | русский | 5 | 2 | 0 | | |
| | татарский | 3 | 3 | 2 | | |
| | английский | 0 | 1 | 10 | | |

Confusion matrix

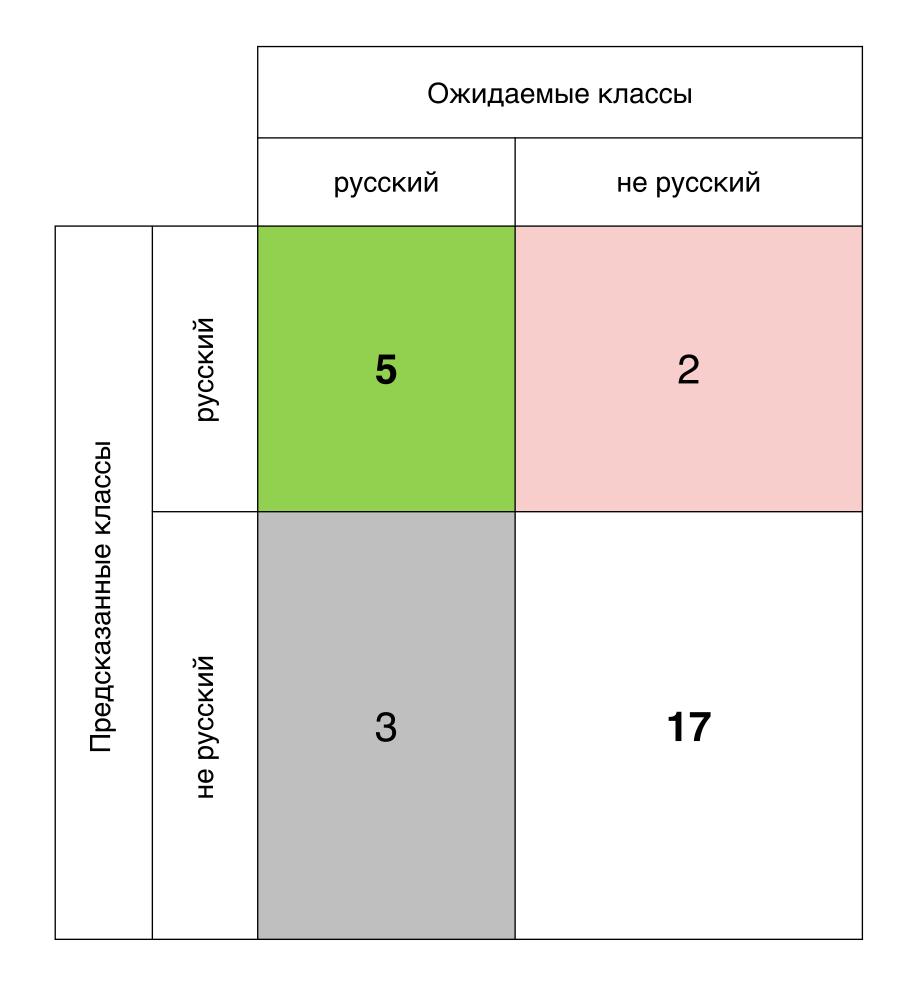
multiclass classification

| | | Ожидаемые классы | | |
|----------------------|------------|------------------|-----------|------------|
| | | русский | татарский | английский |
| Предсказанные классы | русский | 5 | 2 | 0 |
| | татарский | 3 | 3 | 2 |
| | английский | 0 | 1 | 10 |

Precision =
$$\frac{5}{7}$$

$$Recall = \frac{5}{8}$$

$$F_1 = \frac{2}{3}$$



Confusion table

Confusion matrix

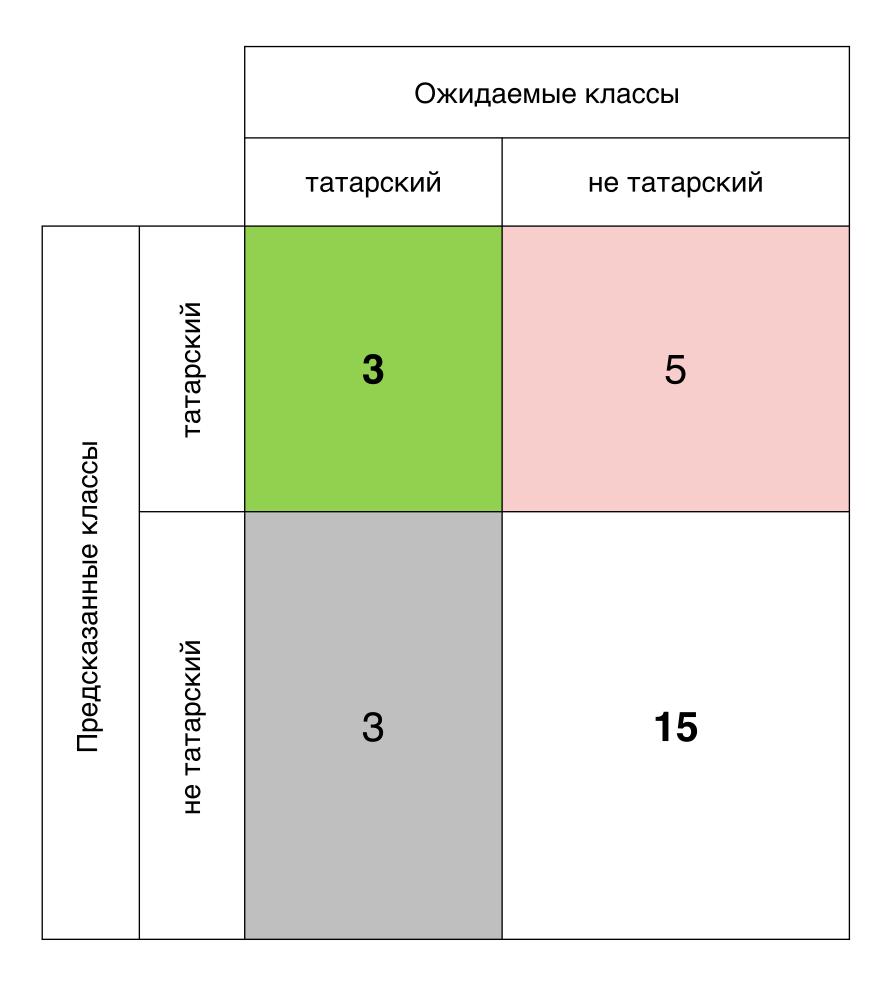
multiclass classification

| | | Ожидаемые классы | | | |
|----------------------|------------|------------------|-----------|------------|--|
| | | русский | татарский | английский | |
| Предсказанные классы | русский | 5 | 2 | 0 | |
| | татарский | 3 | 3 | 2 | |
| | английский | 0 | 1 | 10 | |

$$Precision = \frac{3}{8}$$

Recall =
$$\frac{1}{2}$$

$$F_1 = \frac{3}{7}$$



Confusion table

Confusion matrix

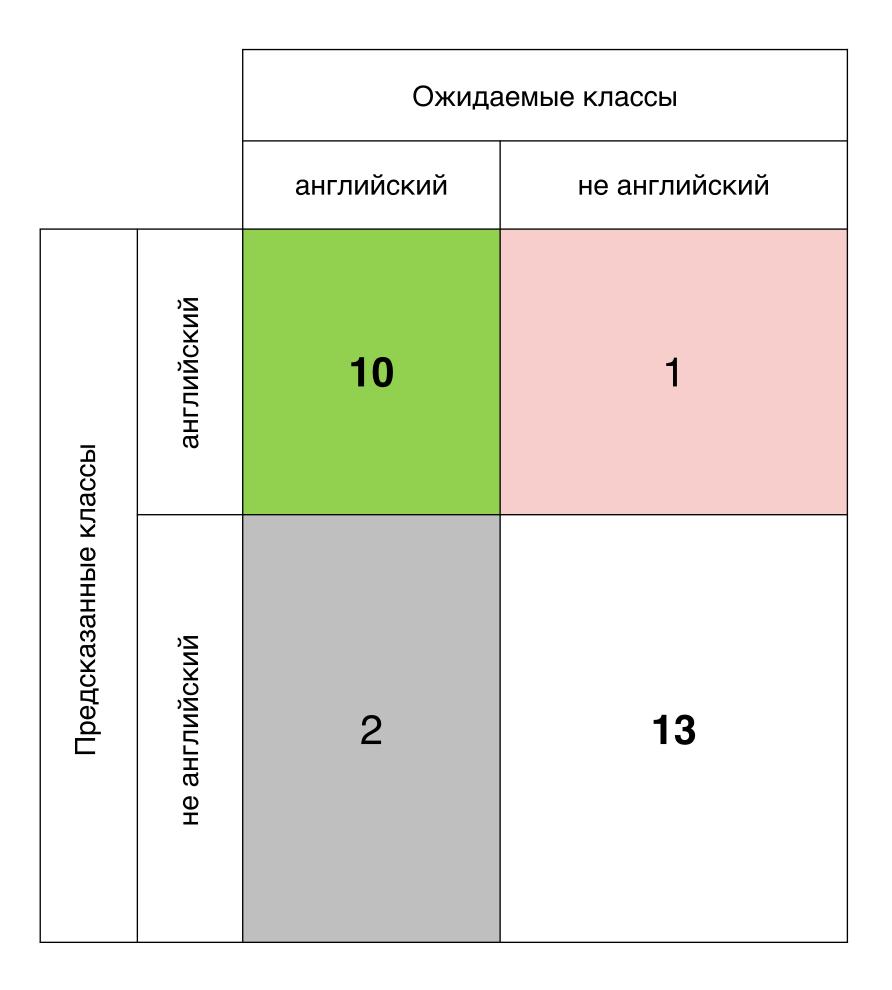
multiclass classification

| | | | ССЫ | |
|----------------------|------------|---------|-----------|------------|
| | | русский | татарский | английский |
| Предсказанные классы | русский | 5 | 2 | 0 |
| | татарский | 3 | 3 | 2 |
| | английский | 0 | 1 | 10 |

$$Precision = \frac{10}{11}$$

Recall =
$$\frac{5}{6}$$

$$F_1 = \frac{20}{23}$$



Confusion table

Confusion matrix

multiclass classification

Макро усреднение

Ожидаемые классы

не русский

Precision =
$$\frac{1}{|C|} \sum \text{Precision}^c$$
 Recall = $\frac{1}{|C|} \sum \text{Recall}^c$ $F_1 = \frac{1}{|C|} \sum F_1^c$

| | | Ожида | емые классы |
|----------------------|--------------|-----------|--------------|
| | | татарский | не татарский |
| классы | татарский | 3 | 5 |
| Предсказанные классы | не татарский | 3 | 15 |

Precision =
$$\frac{1231}{1848} \approx 0.6661$$

$$Recall = \frac{47}{72} \approx 0.6528$$

$$F_1 = \frac{949}{1449} \approx 0.6549$$

| | | Ожида | емые классы |
|----------------------|---------------|------------|---------------|
| | | английский | не английский |
| классы | английский | 10 | 1 |
| Предсказанные классы | не английский | 2 | 13 |

| классы | русский | 5 | 2 |
|----------------------|------------|---|----|
| Предсказанные классы | не русский | 3 | 17 |

русский

multiclass classification

Взвешенное макро усреднение

$$\alpha_c = tp_c + fn_c$$

$$Precision = \frac{\sum \alpha_c \cdot Precision^c}{\sum \alpha_c} \quad Recall = \frac{\sum \alpha_c \cdot Recall^c}{\sum \alpha_c} \quad F_1 = \frac{\sum \alpha_c \cdot F_1^c}{\sum \alpha_c}$$

$$Recall = \frac{\sum \alpha_c \cdot Recall^c}{\sum \alpha_c}$$

$$F_1 = \frac{\sum \alpha_c \cdot F_1^c}{\sum \alpha_c}$$

| | | Ожида | емые классы |
|----------------------|------------|---------|-------------|
| | | русский | не русский |
| классы | русский | 5 | 2 |
| Предсказанные классы | не русский | 3 | 17 |

| | | Ожидаемые классы | |
|----------------------|--------------|------------------|--------------|
| | | татарский | не татарский |
| Классы | татарский | 3 | 5 |
| Предсказанные классы | не татарский | 3 | 15 |

Precision =
$$\frac{5813}{8008} \approx 0.7259$$

$$Recall = \frac{9}{13} \approx 0.6923$$

$$F_1 = \frac{4429}{6279} \approx 0.7054$$

| | | Ожида | емые классы |
|----------------------|---------------|------------|---------------|
| | | английский | не английский |
| слассы | английский | 10 | 1 |
| Предсказанные классы | не английский | 2 | 13 |

Сложности

• Очень короткий текст

| Один два три четыре пять | русский |
|---------------------------|-------------|
| Адзін два тры чатыры пяць | белорусский |
| Один два три чотири п'ять | украинский |
| Един два три четири пет | болгарский |
| Еден два три четири пет | македонский |
| Један два три четири пет | сербский |

Сложности

• Текст сразу на нескольких языках

I heard "El que quiera entender que entienda" by Mägo de Oz?

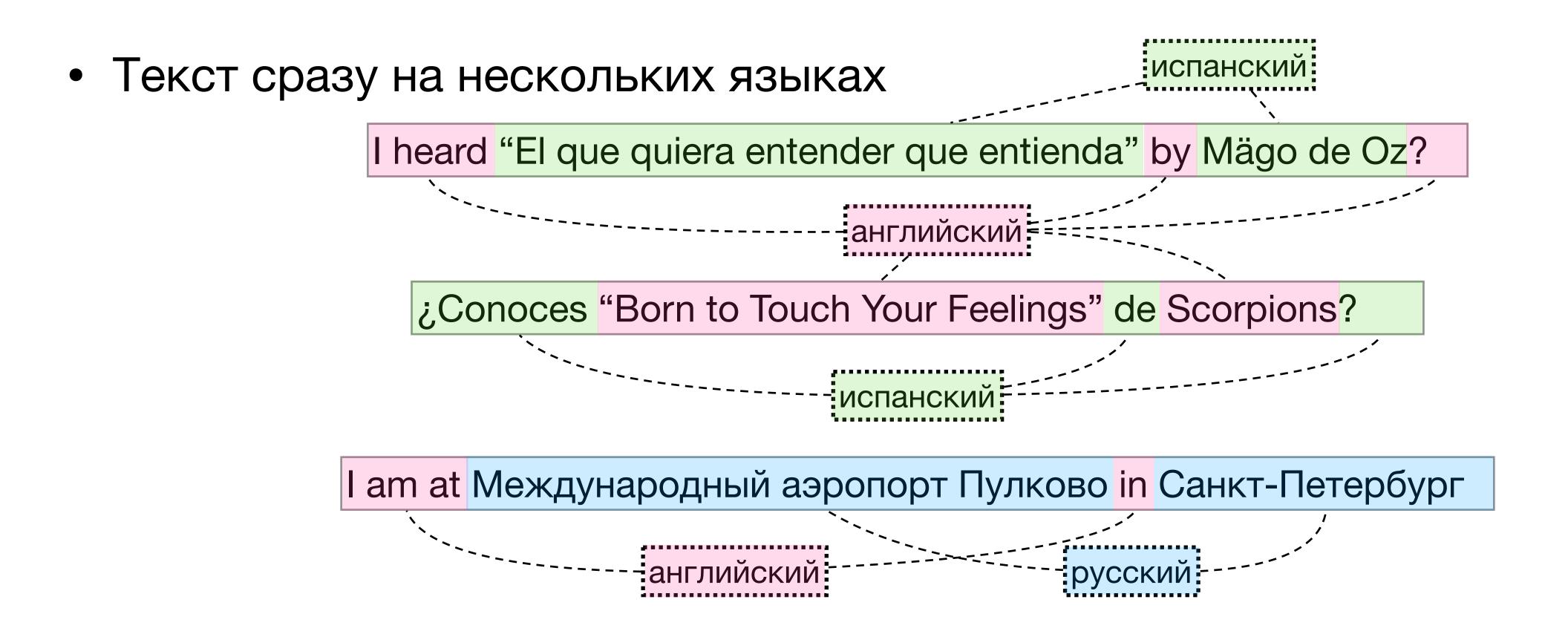
¿Conoces "Born to Touch Your Feelings" de Scorpions?

I am at Международный аэропорт Пулково in Санкт-Петербург

Сложности

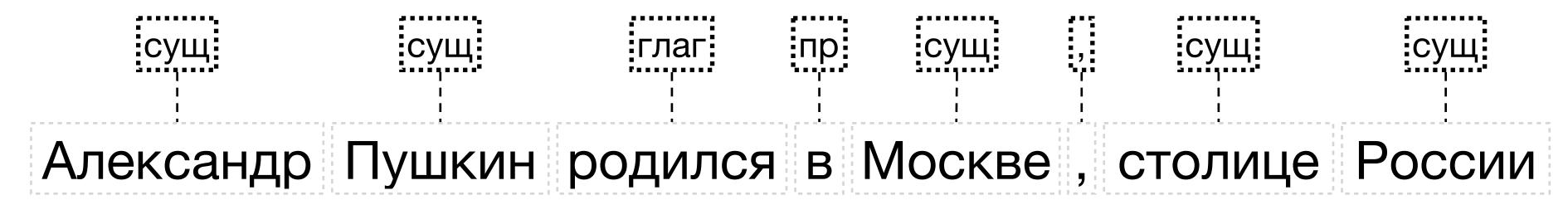
• Текст сразу на нескольких языках

Сложности



- Задача
 - На входе: предложение (последовательность токенов)
 - На выходе: морфологическая метка для каждого слова
- Морфологическая метка часть речи
 - Существительное
 - Прилагательное
 - Глагол

•



Задача разметки последовательности

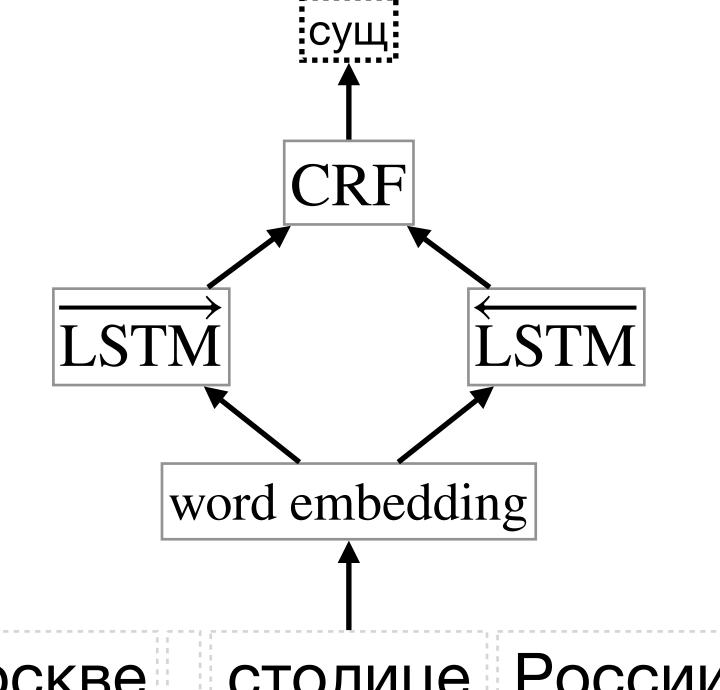
- На входе: последовательность токенов
- На выходе: морфологическая метка для каждого слова

Методы

- CRF
- RNN (biRNN)
- biRNN-CRF

Признаки

- Слова (векторные представления)
- Части слов (префиксы, суффиксы)



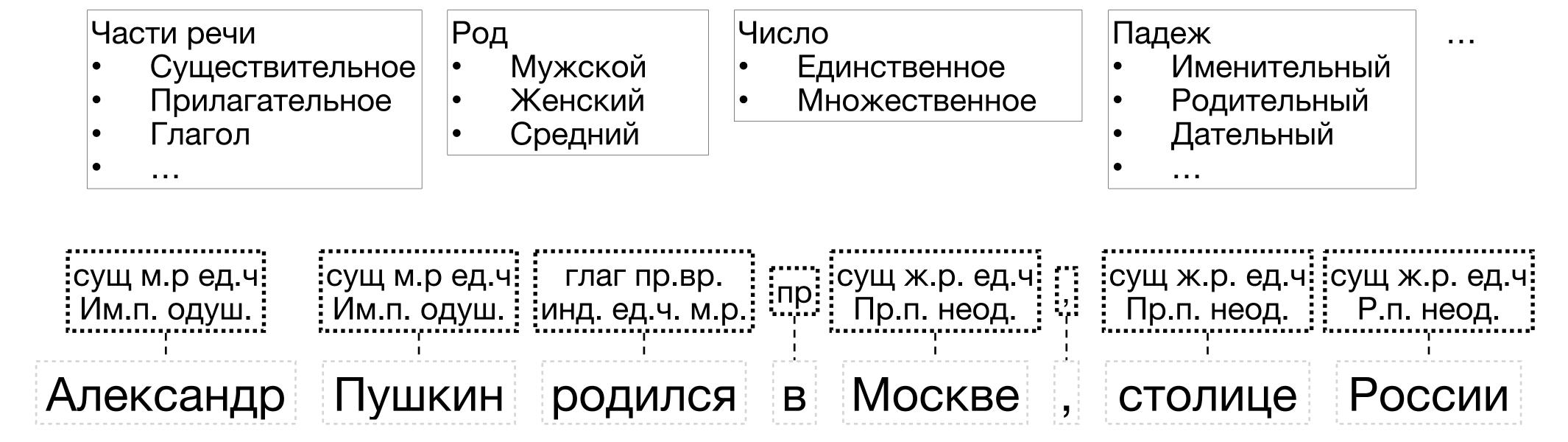
Александр Пушкин родился в Москве, столице России

Методы оценки качества

- Метрики многоклассовой классификации
 - Precision, Recall, F_1

- Sentence Accuracy
 - Количество правильно разобранных предложений среди всех предложений

- Задача
 - На входе: предложение (последовательность токенов)
 - На выходе: морфологическая метка для каждого слова
- Морфологическая метка часть речи + граммемы



Проблемы

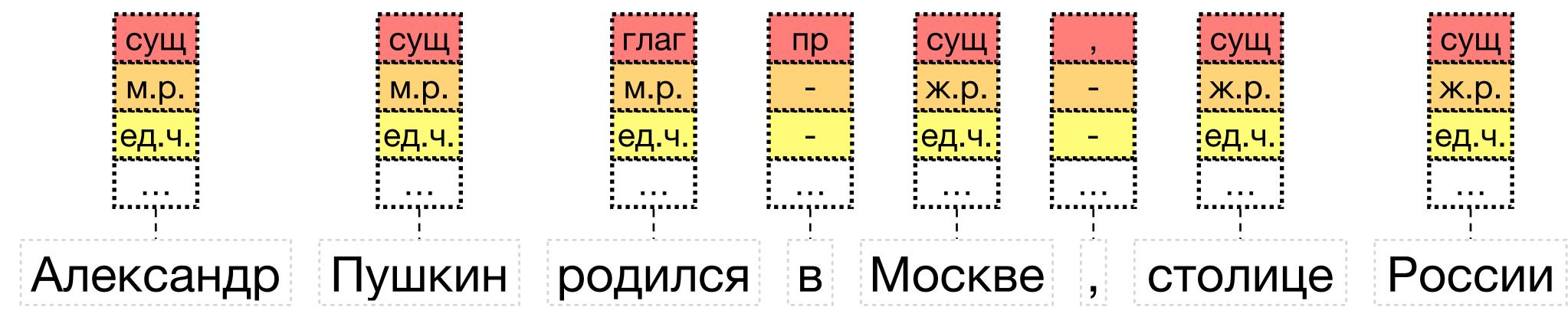
- Много классов (UD SynTagRus: 733)
- Некоторые классы слишком редкие \Rightarrow не получится обучиться

Проблемы

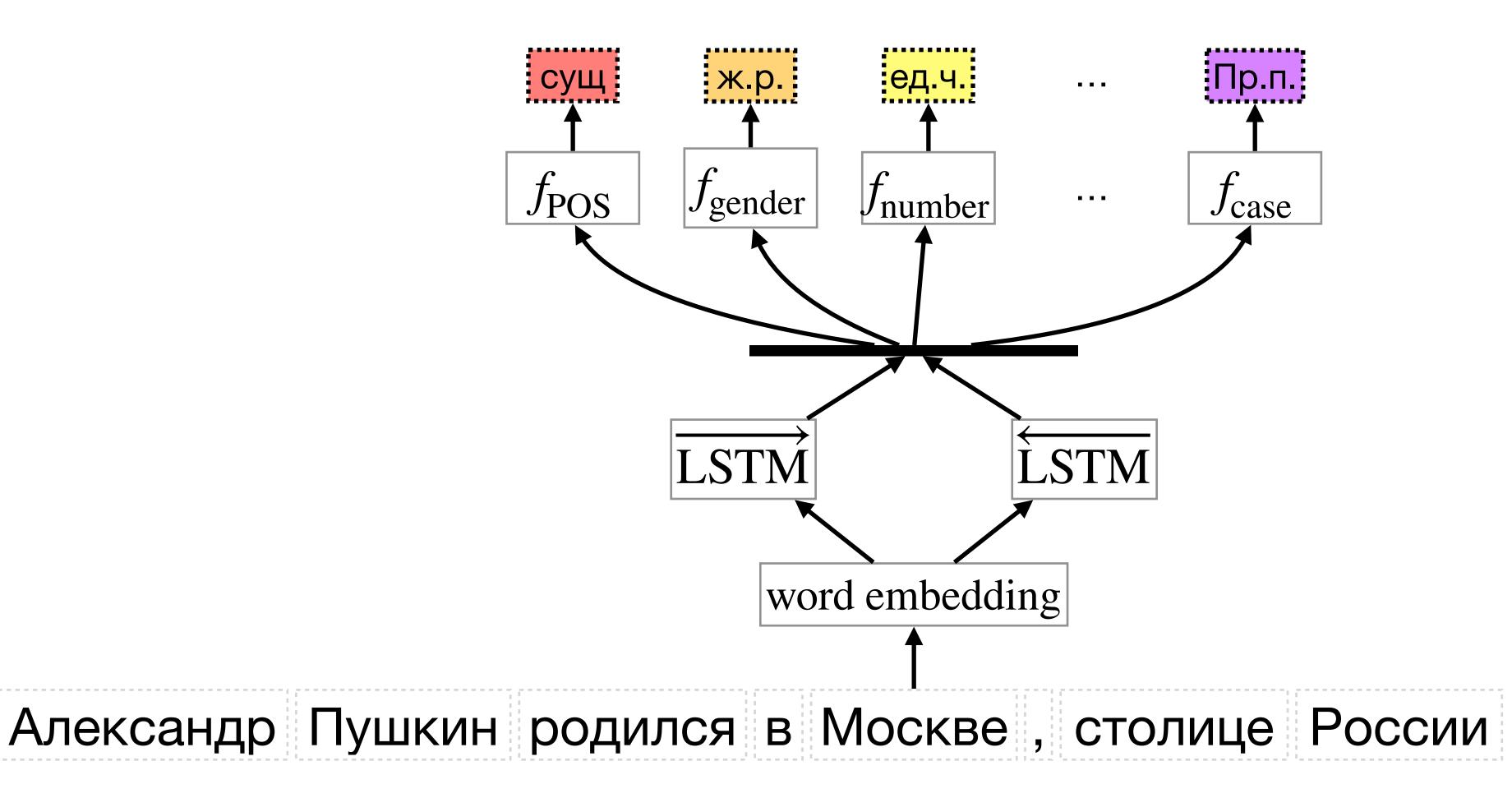
- Много классов (UD SynTagRus: 733)
- Некоторые классы слишком редкие ⇒ не получится обучиться

Решение

- Разделить метки по грамматическим категориям
- Назначить каждому токену метку из каждой категории
- Multilabel classification



Multilabel classification



Multilabel classification

При таком подходе может получиться неправильная метка:

- Род, число падеж у предлогов/союзов
- Род у глаголов настоящего времени
- Время у существительных

Решение:

- Собрать множество возможных меток (из корпуса/руками)
- Определять наиболее вероятную метку из возможных

$$p$$
(сущ, м.р, ..., Пр.п | w) = p (сущ | w) $\cdot p$ (м.р | w) $\cdot ... \cdot p$ (Пр.п | w)

Задача

- На входе: предложение/последовательность слов (токенов)
- На выходе: последовательность лемм слов (токенов) в нормальной форме

Нормальная форма слова – каноническая форма слова

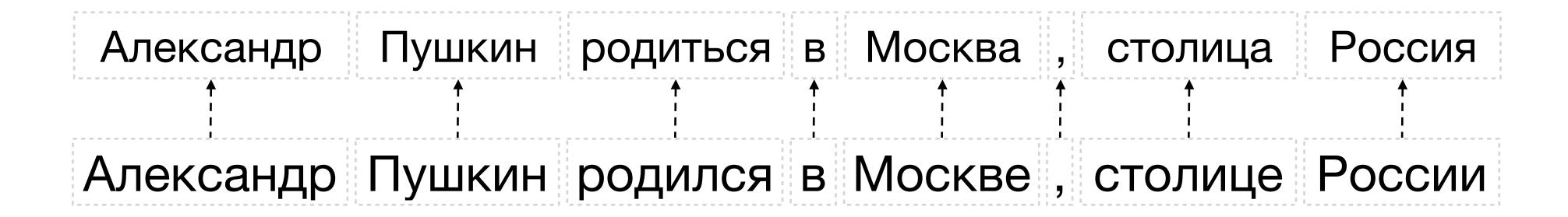
- Для существительных:
 - единственное число
 - именительный падеж
- Для глаголов:
 - инфинитив

• . . .

Задача

- На входе: предложение/последовательность слов (токенов)
- На выходе: последовательность лемм слов (токенов) в нормальной форме

Задача похожа на разметку последовательности

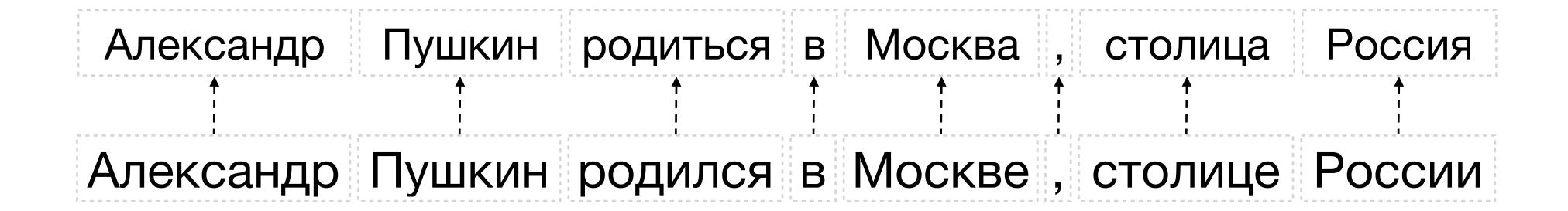


Задача

- На входе: предложение/последовательность слов (токенов)
- На выходе: последовательность лемм слов (токенов) в нормальной форме

Задача похожа на разметку последовательности, но

- На выходе слова, а не классы
- Слова на выходе не связаны между собой

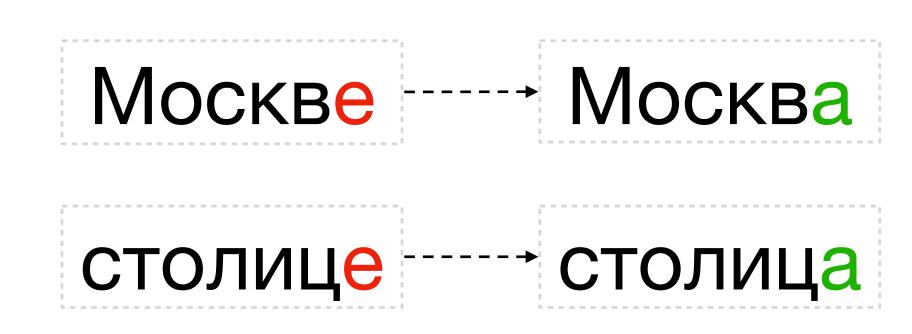


Задача

- На входе: слово (токен)
- На выходе: лемма слово (токен) в нормальной форме

Процесс построения леммы

- Удалить суффикс словоформы
- Добавить суффикс леммы
- Удалить префикс словоформы
- Добавить префикс леммы



Методы, основанные на правилах

• LemmaGen (2010)

IF suffix THEN transformation EXCEPT exceptions

```
IF "" THEN ""→" EXCEPT
IF "d" THEN "d"→""
ELSE IF "ote" THEN "ote"→"ite"
ELSE IF "ing" THEN "ing"→"e" EXCEPT
IF "ting" THEN ""→""
ELSE IF "ten" THEN "ten"→"e"
ELSE IF "s" THEN "s"→""
```

Многоклассовая классификация

Классы:

- Суффиксы словоформ, суффиксы лемм, префиксы словоформ, префиксы лемм
- Правила преобразования суффикса, правила преобразования префикса

Признаки:

- Символы слова
- Морфологические признаки

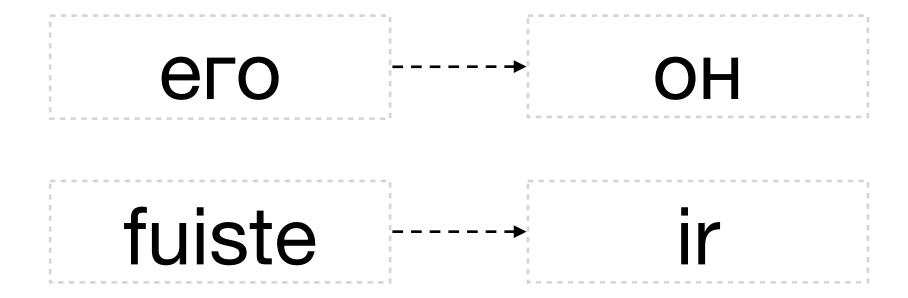
•

Классификаторы:

- CNN
- Perceptron

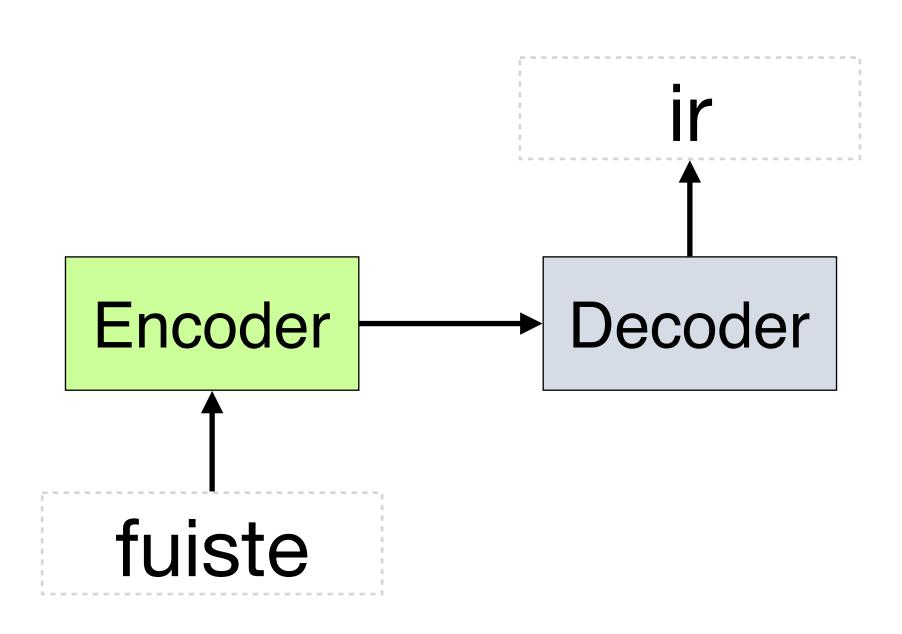
•

• Важная проблема таких подходов — слова исключения

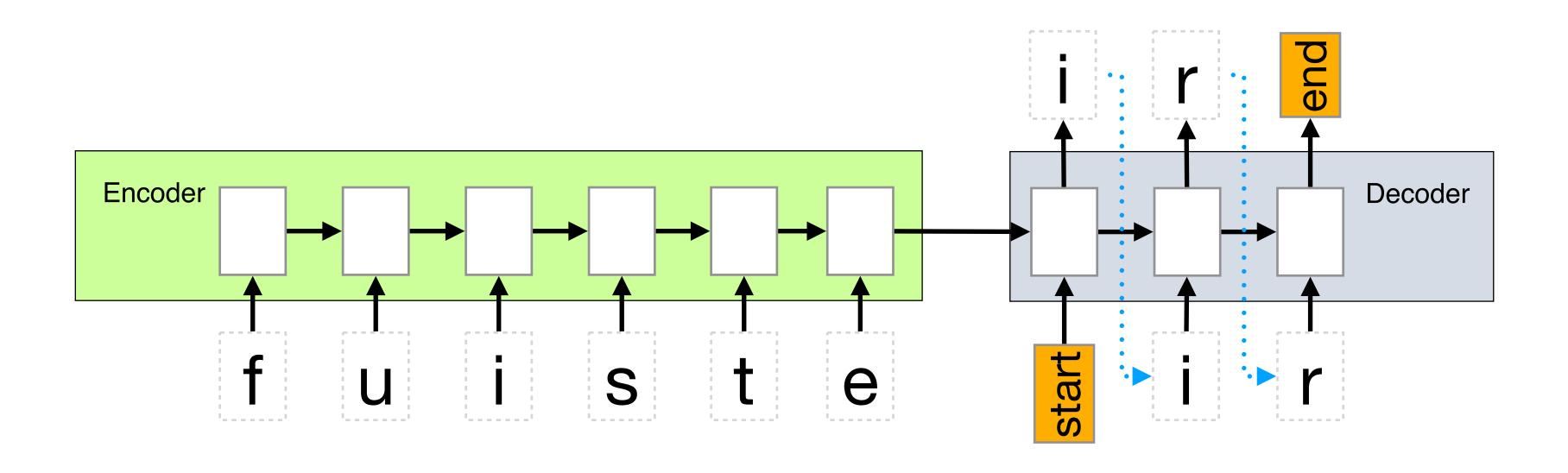


• Решение — словари исключений

Нормализация (лемматизация) слов seq2seq



Нормализация (лемматизация) слов seq2seq



Многозначность

Грамматическая омонимия

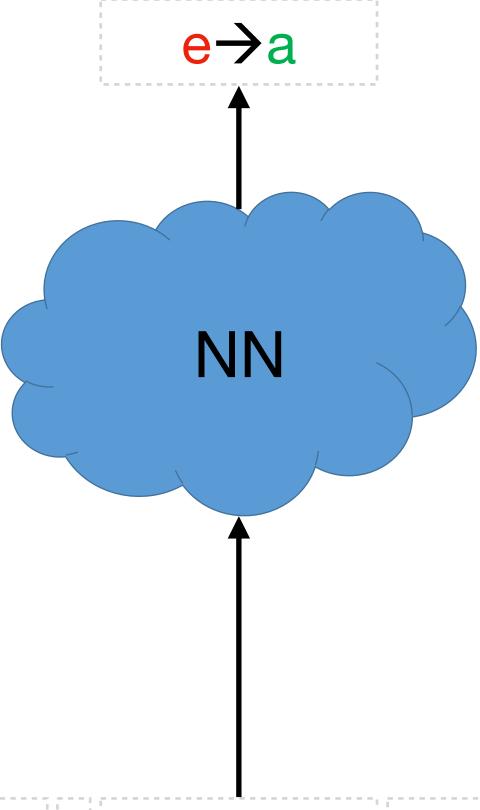


Решение

• Использовать информацию о морфологии слова

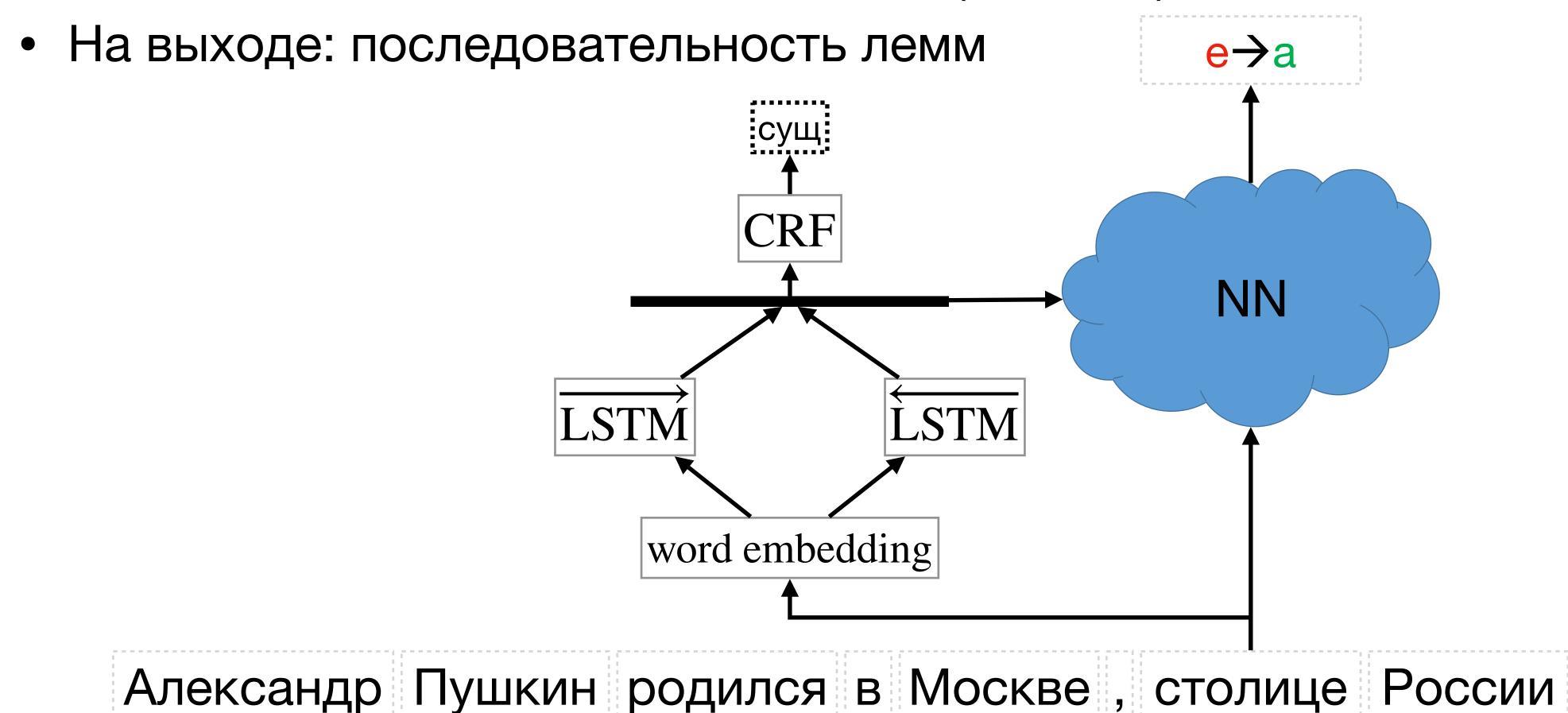
Задача

- На входе: последовательность слов (токенов)
- На выходе: последовательность лемм



Задача

• На входе: последовательность слов (токенов)



Оценка качества

• Классическая точность (Accuracy)

$$Accuracy = \frac{correct}{total}$$

• Точность среди неизвестных слов (out-of-vocabulary)

$$Accuracy_{OOV} = \frac{correct_{OOV}}{total_{OOV}}$$

Следующая лекция

Синтаксический анализ

- Дерево составляющих
 - Алгоритм СКҮ
- Дерево зависимостей
 - transition-based parsing
 - графовые методы