

ЯЗЫКОВЫЕ МОДЕЛИ

КВАНТИТАТИВНЫЙ АНАЛИЗ ТЕКСТА

Кирилл Александрович Маслинский

16.05.2022 / 09

ЕУСП6

КОНТЕКСТ

КОНТЕКСТ

МОДЕЛЬ КОНТЕКСТА: N-ГРАММЫ

N последовательно стоящих друг за другом слов.

Униграммы

Восторг внезапный ум пленил .

Биграммы

Восторг внезапный ум пленил <.>

Триграммы

<s> Восторг внезапный ум пленил <.>

N последовательно стоящих друг за другом слов.

Униграммы

Восторг внезапный ум пленил .

Биграммы

Восторг внезапный ум пленил <.>

Триграммы

<s> Восторг внезапный ум пленил <.>

N последовательно стоящих друг за другом слов.

Униграммы

Восторг **внезапный** ум пленил .

Биграммы

Восторг **внезапный** ум пленил <.>

Триграммы

<s> Восторг **внезапный** ум пленил <.>

N последовательно стоящих друг за другом слов.

Униграммы

Восторг внезапный ум пленил .

Биграммы

Восторг внезапный ум пленил <.>

Триграммы

<s> Восторг внезапный ум пленил <.>

N последовательно стоящих друг за другом слов.

Униграммы

Восторг внезапный ум **п**ленил .

Биграммы

Восторг внезапный ум **п**ленил <.>

Триграммы

<s> Восторг внезапный **ум** пленил <.>

ВЕРОЯТНОСТЬ ЯЗЫКОВЫХ СОБЫТИЙ

искусственный ...?

искусственный интеллект?

а в корпусе русской детской литературы XX—XXI вв.?

ВЕРОЯТНОСТЬ ЯЗЫКОВЫХ СОБЫТИЙ

- Классическая вероятность: монетка или кость
- Вероятность, основанная на частотности (урна с шарами, мешок слов)
- В лингвистике считаем события в корпусе
- вероятность = относительная частотность

Пример расчета вероятности слова

Всего слов в корпусе = 46,804,371

искусственный = 121

$$P(\text{искусственный}) = \frac{121}{46804371} \approx 0.000002585 = 2.585 \text{ IPM}$$

ВЕРОЯТНОСТЬ ЯЗЫКОВЫХ СОБЫТИЙ

- Классическая вероятность: монетка или кость
- Вероятность, основанная на частотности (урна с шарами, мешок слов)
- В лингвистике считаем события в корпусе
- вероятность = относительная частотность

Пример расчета вероятности слова

Всего слов в корпусе = 46,804,371

искусственный = 121

$$P(\text{искусственный}) = \frac{121}{46804371} \approx 0.000002585 = 2.585 \text{ IPM}$$

- $P(\text{искусственный}) = \frac{121}{46804371}$
- $P(\text{интеллект}) = \frac{179}{46804371}$
- $P(\text{искусственный интеллект}) = ?$

$$P(A \wedge B) = P(A) \cdot P(B) \quad | \quad \text{IFF } A \text{ и } B \text{ независимы} \quad (1)$$

$$P(A \wedge B) = P(A) \cdot P(B) \quad | \quad \text{IFF } A \text{ и } B \text{ независимы} \quad (1)$$

$$\begin{aligned} P(\text{искусственный интеллект}) &= \frac{121}{46804371} \cdot \frac{179}{46804371} = \\ &= 0.000000000000989 = 0.00000989 \text{ IPM} \end{aligned}$$

$$P(B|A) = \frac{P(B \wedge A)}{P(A)} \quad (2)$$

$$P(\text{интеллект}|\text{искусственный}) = \frac{P(\text{искусственный интеллект})}{P(\text{искусственный})} =$$

$$= \frac{\frac{4}{46804371}}{\frac{121}{46804371}} = \frac{4}{121} = 0.033$$

сильный искусственный ...?

сильный искусственный интеллект!

ПРЕДСКАЗАНИЕ СЛОВА В КОНТЕКСТЕ

сильный искусственный

$$P(\text{интеллект} | \text{сильный искусственный}) =$$

$$\frac{f(\text{сильный искусственный интеллект})}{f(\text{сильный искусственный})} = \frac{?}{?} = ?$$

ВЕРОЯТНОСТЬ ЯЗЫКОВЫХ СОБЫТИЙ

Коллокации

Модель для коллокаций

Слова, встречающиеся в текстах рядом чаще, чем можно ожидать в результате случайности.

	интеллект	\neg интеллект
искусственный		
\neg искусственный		

ВЕРОЯТНОСТЬ ЯЗЫКОВЫХ СОБЫТИЙ

ЯЗЫКОВАЯ МОДЕЛЬ

- $P(\text{Создан сильный искусственный интеллект.}) = ?$
- $P(\text{Создан сильный искусственный крокодил.}) = ?$

Языковая модель — приписывает вероятность фрагменту текста (высказыванию, предложению...)

В хорошей модели вероятности языковых фрагментов соответствуют их относительной частотности в текстах.

Иными словами:

- максимизирует вероятность реальных текстов
- минимизирует вероятность нереальных текстов

УНИГРАММНАЯ ЯЗЫКОВАЯ МОДЕЛЬ

Создан³¹ сильный¹⁹⁰⁵ искусственный¹²¹ интеллект⁷¹!

Создан³¹ сильный¹⁹⁰⁵ искусственный¹²¹ крокодил⁴⁴²!

УНИГРАММНАЯ ЯЗЫКОВАЯ МОДЕЛЬ

Создан³¹ сильный¹⁹⁰⁵ искусственный¹²¹ интеллект⁷¹!

Создан³¹ сильный¹⁹⁰⁵ искусственный¹²¹ крокодил⁴⁴²!

$$\frac{31}{46804371} \cdot \frac{1905}{46804371} \cdot \frac{121}{46804371} \cdot \frac{71}{46804371} = 1.06 \times 10^{-22}$$

$$\frac{31}{46804371} \cdot \frac{1905}{46804371} \cdot \frac{121}{46804371} \cdot \frac{442}{46804371} = 6.58 \times 10^{-22}$$

Крокодил в шесть раз вероятнее!

$$P(A, B) = P(B|A)P(A) \quad (3)$$

$$P(W_1, W_2, W_3, W_4) =$$

$$P(W_4|W_1, W_2, W_3) \cdot P(W_3|W_1, W_2) \cdot P(W_2|W_1) \cdot P(W_1)$$

$P(\text{Создан сильный искусственный крокодил}) = P(\text{Создан})$
 $P(\text{сильный}|\text{Создан}) P(\text{искусственный}|\text{Создан сильный})$
 $P(\text{крокодил}|\text{Создан сильный искусственный})$

- система с конечным числом состояний
- следующее состояние зависит только от текущего

Применительно к тексту:

Следующее слово зависит только от предыдущего (N предыдущих)

Markov assumption:

$$P(\text{Создан сильный искусственный интеллект}) \approx$$

$$P(\text{интеллект}|\text{искусственный}) \cdot P(\text{искусственный}|\text{сильный})$$

$$\cdot P(\text{сильный}|\text{создан}) =$$

$$\frac{2.21}{45.65} \cdot \frac{0.01}{189.7} \cdot \frac{0.003}{16.28} \approx 4.7 \times 10^{-10}$$

Markov assumption:

$$P(\text{Создан сильный искусственный крокодил}) \approx$$

$$P(\text{крокодил}|\text{искусственный}) \cdot P(\text{искусственный}|\text{сильный})$$

$$\cdot P(\text{сильный}|\text{создан}) =$$

$$\frac{0}{45.65} \cdot \frac{0.01}{189.7} \cdot \frac{0.003}{16.28} \approx 0$$

Markov assumption:

$$P(\text{Создан сильный искусственный крокодил}) \approx$$

$$P(\text{крокодил}|\text{искусственный}) \cdot P(\text{искусственный}|\text{сильный})$$

$$\cdot P(\text{сильный}|\text{создан}) =$$

$$\frac{0.001}{45.65} \cdot \frac{0.01}{189.7} \cdot \frac{0.003}{16.28} \approx 2.13 \times 10^{-13}$$

Проблема недостаточных данных (sparse data)

- В результате обучения на конечном корпусе очень многие N-граммы получают нулевую вероятность
- Хотя в действительности должны иметь ненулевую (встречаются в текстах)

Если в предложении встречается слово, отсутствовавшее в корпусе, вероятность такого предложения в модели **равна нулю!**

Задача

Добиться, чтобы модель приписывала ненулевую вероятность любому тексту.

Добавить 1 ко всем частотам в корпусе

Для униграммной модели (сглаженная относительная частотность слова):

$$P_{Laplace} = \frac{f + 1}{N + V} \quad (4)$$

Идея

Оценить вероятность никогда не встречавшихся N-грамм на основании вероятности N-грамм, встречавшихся один раз

- Good-Turing discounting
- Kneser-Ney smoothing

Markov assumption: $P(\text{Создан сильный искусственный крокодил})$
 $\approx P(\text{крокодил}|\text{искусственный}) P(\text{искусственный}|\text{сильный})$
 $P(\text{сильный}|\text{создан})$

$$\frac{0.001}{45.65} \cdot \frac{0.01}{189.7} \cdot \frac{0.003}{16.28} \approx 2.13 \times 10^{-13}$$

$$\frac{0.001}{45.2} \cdot \frac{0.01}{187.8} \cdot \frac{0.004}{16.12} \approx 2.92 \times 10^{-13}$$

Интеллект:

$$4.7 \times 10^{-10}$$

$$6.4 \times 10^{-10}$$

$P(\text{Создан сильный искусственный интеллект}) \approx$

$P(\text{интеллект} | \text{сильный искусственный}) \cdot$

$P(\text{искусственный} | \text{создан сильный})$

Конечный автомат

