

word2vec

ПОХОЖИЕ СЛОВА

- «Идти» и «шагать» — синонимы
- Для компьютера это разные строки
- Как понять, что они похожи?

- › «Идти» и «шагать» — синонимы
- › Для компьютера это разные строки
- › Как понять, что они похожи?
- › На основе данных!
- › Слова со схожим смыслом часто идут в паре с одними и теми же словами
- › У них похожие контексты

- › Хотим каждое слово представить как вещественный вектор: $w \rightarrow \vec{w} \in \mathbb{R}^d$
- › Какие требования?
 - ▶ Размерность d должна быть не очень велика
 - ▶ Похожие слова должны иметь близкие векторы
 - ▶ Арифметические операции над векторами должны иметь смысл

- › Будем обучать представления слов так, чтобы они хорошо предсказывали свой контекст
- › Выборка состоит из текстов, каждый представляет собой последовательность слов $w_1, \dots, w_i, \dots, w_n$

$$\sum_{i=1}^n \sum_{j=-k}^k \log p(w_{i+j} | w_i) \rightarrow \max,$$

- › Выборка состоит из текстов, каждый представляет собой последовательность слов $w_1, \dots, w_i, \dots, w_n$

$$\sum_{i=1}^n \sum_{j=-k}^k \log p(w_{i+j}|w_i) \rightarrow \max,$$

где вероятность вычисляется через soft-max:

$$p(w_i|w_j) = \frac{\exp(\langle \vec{w}_i, \vec{w}_j \rangle)}{\sum_w \exp(\langle \vec{w}, \vec{w}_j \rangle)}$$

- › Косинусное расстояние хорошо отражает схожесть слов по тематике (в зависимости от корпуса)
- › $\vec{\text{king}} - \vec{\text{man}} + \vec{\text{woman}} \approx \vec{\text{queen}}$
- › $\vec{\text{Moscow}} - \vec{\text{Russia}} + \vec{\text{England}} \approx \vec{\text{London}}$
- › Перевод: $\vec{\text{oñe}} - \vec{\text{uño}} + \vec{\text{four}} \approx \vec{\text{quarto}}$
- › Среднее представление по всем словам в тексте — хорошее признаковое описание

- › Проблема мешка слов — слишком большое количество признаков
- › Средний word2vec-вектор позволяет получить компактное признаковое описание
- › При размерности вектора 100 можно обучать композиции деревьев

- › word2vec позволяет описать каждое слово вектором
- › Похожие слова имеют близкие векторы
- › Признаки для текста — средний вектор по всем словам