

Анализ текстов

Лекция 11



Как научить компьютер читать?



Разрешение анафоры

“Мама вымыла раму, и теперь она блестит”

“Мама вымыла раму, и теперь она устала”



Задачи обработки текстов

1. Синтаксический анализ
2. Семантический анализ (классификация)
3. Порождение текста



Синтаксический анализ

1. Разметка частей речи
2. Морфологическая сегментация
3. Выделение границ предложений
4. Разрешение кореференций



Семантический анализ

1. Предсказание следующего слова по контексту
2. Анализ тональности текста
3. Ответы на вопросы (когда?, где?)
4. Выделение фактов
5. Поиск синонимов



Порождение текста

1. Автореферирование
2. Машинный перевод
3. Диалоговые модели



Классические модели обработки текста



Мешок слов

“Сантехник чинит кран”, “Сантехник из водоканала”,
“Кран течёт”, “Электрик Анатолий”, “Анатолий чинит
лампу”.

сантехник	чинить	кран	водоканал	течёт	электрик	лампа	Анатолий
1	1	1					
1			1				
		1		1			
					1		1
	1					1	1



TF-IDF

$$\text{tf}(t, d) = \frac{n_t}{\sum_k n_k}, \quad \text{idf}(t, D) = \log \frac{|D|}{|\{d_i \in D \mid t \in d_i\}|}$$

$$\text{tf-idf}(t, d, D) = \text{tf}(t, d) \times \text{idf}(t, D)$$

<https://ru.wikipedia.org/wiki/TF-IDF>



Учёт порядка слов

“Мама мыла раму с мылом”

- n-граммы слов (“мама мыла”, “мыла раму”)
- k-skip-n-граммы (“мама раму”, “мыла мылом”)
- Буквенные n-граммы (“мыл”, “рам”, “мыл”)



Проблемы классической модели

- Вектора признаков очень разреженные
- Много 0 в признаках
- Количество признаков = размер словаря \sim миллионы
- В основе - гипотеза независимости слов

Итог: векторы слов не отражают семантику слов



Векторные представления

Word Embedding



Дистрибутивная гипотеза

“You shall know a word by the company it keeps” – Firth, 1957.

Лингвистические единицы, встречающиеся в схожих контекстах, имеют близкие значения. (1960-ые)

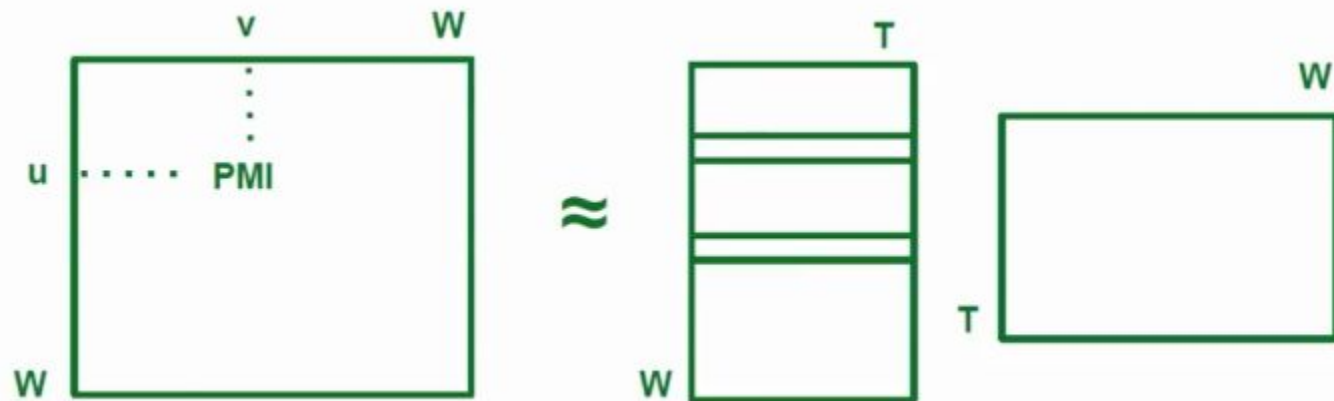
government debt problems turning into banking crises as has happened in
saying that Europe needs unified banking regulation to replace the hodgepodge

↖ These words will represent *banking* ↗

<https://cs224d.stanford.edu/lectures/CS224d-Lecture2.pdf>



Латентный семантический анализ LSA



Сингулярное разложение матрицы

$$\begin{array}{ccccc}
 \begin{array}{c} m \\ \boxed{} \\ n \\ X \end{array} & = & \begin{array}{c} r \\ \boxed{\begin{array}{c} | \quad | \quad | \quad \dots \\ U_1 \ U_2 \ U_3 \ \dots \\ | \quad | \quad | \end{array}} \\ n \quad U \end{array} & \begin{array}{c} r \\ \boxed{\begin{array}{ccc} S_1 & & \\ & S_2 & \\ & & 0 \end{array}} \\ r \quad S \end{array} & \begin{array}{c} m \\ \boxed{\begin{array}{c} \text{---} V_1 \text{---} \\ \text{---} V_2 \text{---} \\ \text{---} V_3 \text{---} \\ \vdots \end{array}} \\ r \quad V^T \end{array} \\
 \\
 \begin{array}{c} m \\ \boxed{\phantom{\hat{X}}} \\ n \\ \hat{X} \end{array} & = & \begin{array}{c} k \\ \boxed{\begin{array}{c} | \quad | \quad | \quad \dots \\ U_1 \ U_2 \ U_3 \ \dots \\ | \quad | \quad | \end{array}} \\ n \quad \hat{U} \end{array} & \begin{array}{c} k \\ \boxed{\begin{array}{ccc} S_1 & & \\ & S_2 & \\ & & 0 \end{array}} \\ k \quad \hat{S} \end{array} & \begin{array}{c} m \\ \boxed{\begin{array}{c} \text{---} V_1 \text{---} \\ \text{---} V_2 \text{---} \\ \text{---} V_3 \text{---} \\ \vdots \end{array}} \\ k \quad \hat{V}^T \end{array}
 \end{array}$$

<https://cs224d.stanford.edu/lectures/CS224d-Lecture2.pdf>



Распределенное представление

One-Hot encoding

Сантехник = $[1, 0, 0, 0, 0, \dots, 1, \dots, 0]$

Электрик = $[0, 0, 1, 0, 0, \dots, 1, \dots, 1]$

Word Embedding

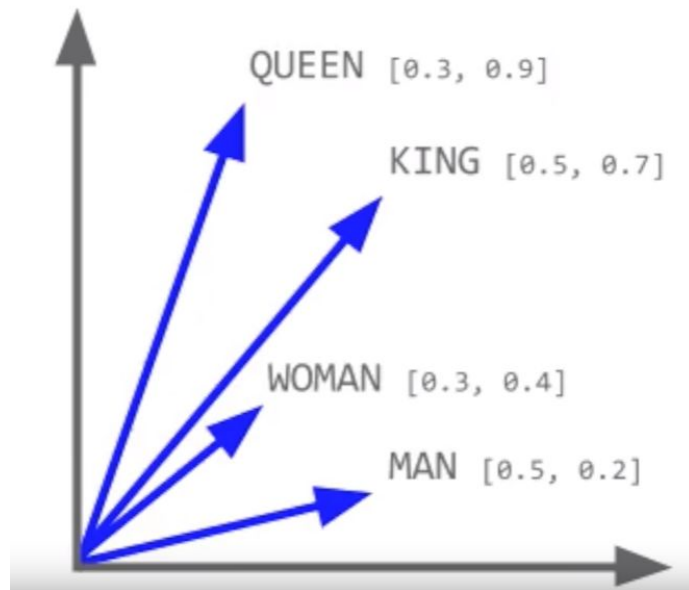
Кран = $[0.22, 0.42, 0.1, 0.5]$

Труба = $[0.23, 0.4, 0.15, 0.7]$

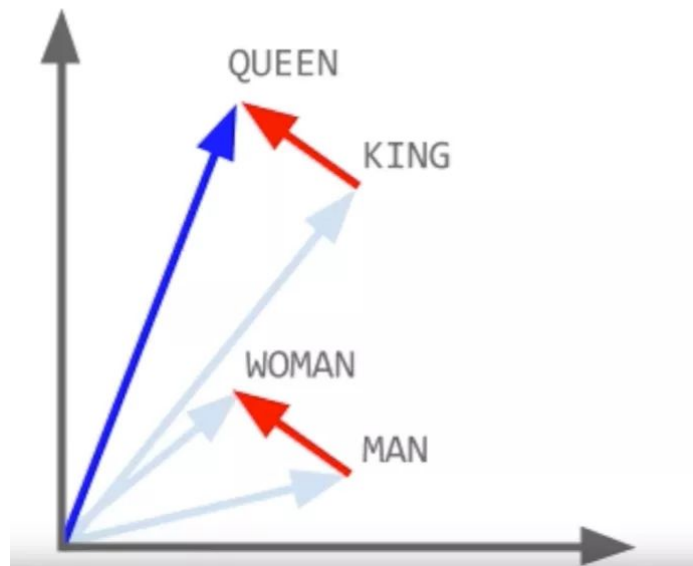


word2vec

Load up the word vectors



So $\text{king} + \text{man} - \text{woman} = \text{queen!}$



word2vec алгоритм

- Каждому слову v сопоставляем вектор \mathbf{w} размера d
- Считаем вероятности $p(v|c_1...c_n)$, где c_i - контекст
- $p(v|c_1...c_n) = f(\mathbf{w}, \mathbf{o})$, \mathbf{o} - параметры функции
- Как найти \mathbf{o} и \mathbf{w} - нейросеть его знает...



Word2vec

...an efficient method for learning high quality distributed vector ...

Context focus word Context

<https://blog.acolyer.org/2016/04/21/the-amazing-power-of-word-vectors/>



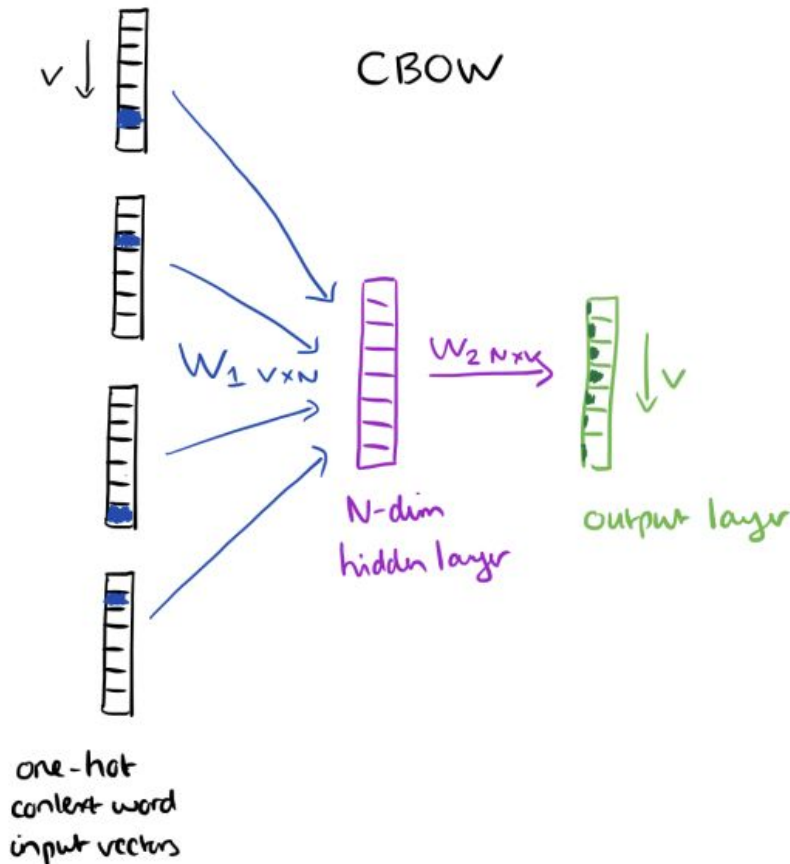
word2vec: CBOW

непрерывный мешок слов

Вход: контекст

Выход: целевое слово

Скрытый слой: векторное представление слов



<https://blog.acolyer.org/2016/04/21/the-amazing-power-of-word-vectors/>



Модификации

- word2vec
 - CBOW
 - skip-gram
 - GLOVE - <https://nlp.stanford.edu/projects/glove/>
- doc2vec
- ...



Литература

1. Глубокое обучение. С. Николенко, А. Кадурын, Е. Архангельская.
2. Векторные представления слов в документах:
<https://www.youtube.com/watch?v=hiDBnEyoZS4>
3. Мастер класс по Gensim: <https://www.youtube.com/watch?v=oBb9aFmp0Hs>
4. Лекции 4 семестра специализации ML&DS:
<https://www.coursera.org/specializations/machine-learning-data-analysis>
5. <https://cs224d.stanford.edu/lectures/CS224d-Lecture2.pdf>
6. <https://github.com/RaRe-Technologies/movie-plots-by-genre>

