



ОБРАЗОВАТЕЛЬНЫЙ ИНТЕНСИВ

Остров 10-22

10-22 июля 2019

Векторные

представления слов.

Классификация

ТЕКСТОВ

Баймурзина Диляра

13.07.2019



Векторные представления слов

Постановка задачи векторизации текста

- Обработка естественного языка (Natural language processing, NLP) предполагает представление текста в некотором машиночитаемом формате.
- Слово традиционно представляется в качестве основной единицы языка.
- Как можно представить слово в численном виде?

One-hot вектора слова

- Составим словарь возможных слов V . Для i -го слова w_i в словаре вектор слова будет иметь вид $v_i = \{0, \dots, 0_{i-1}, 1_i, 0_{i+1}, \dots, 0\}$.
- Размерность вектора равна размеру словаря (например, 200 000).
- Для каждого слова единственная компонента вектора является единицей, все остальные элементы вектора являются нулями

Проблемы:

- Вектора большой размерности замедляют обработку текстов.
- Расстояние между любыми двумя словами одинаковое, не зависит от семантики¹ слов.

motel = [0 0 0 0 0 0 0 0 0 0 1 0 0 0 0]

hotel = [0 0 0 0 0 0 0 1 0 0 0 0 0 0 0]

¹Семантика - раздел лингвистики, изучающий смысловое значение единиц языка.

TF-IDF коэффициенты для векторного представления документов

- Дан набор документов D
- Term frequency — насколько слово w важно для документа d :

$$\text{TF}(w, d) = |\{w \in d\}|$$

- Inverse document frequency — насколько слово w специфично для документа d :

$$\text{IDF}(w, D) = -\log P(w | D) = \log \frac{|D|}{|\{d \in D : w \in d\}|}$$

- TF-IDF для слова w в документе d в наборе текстов D :

$$\text{TF-IDF}(w, d, D) = \text{TF} \cdot \text{IDF}$$

TF-IDF коэффициенты для векторного представления документов

Таким образом получаем вектор документа d в наборе документов D :

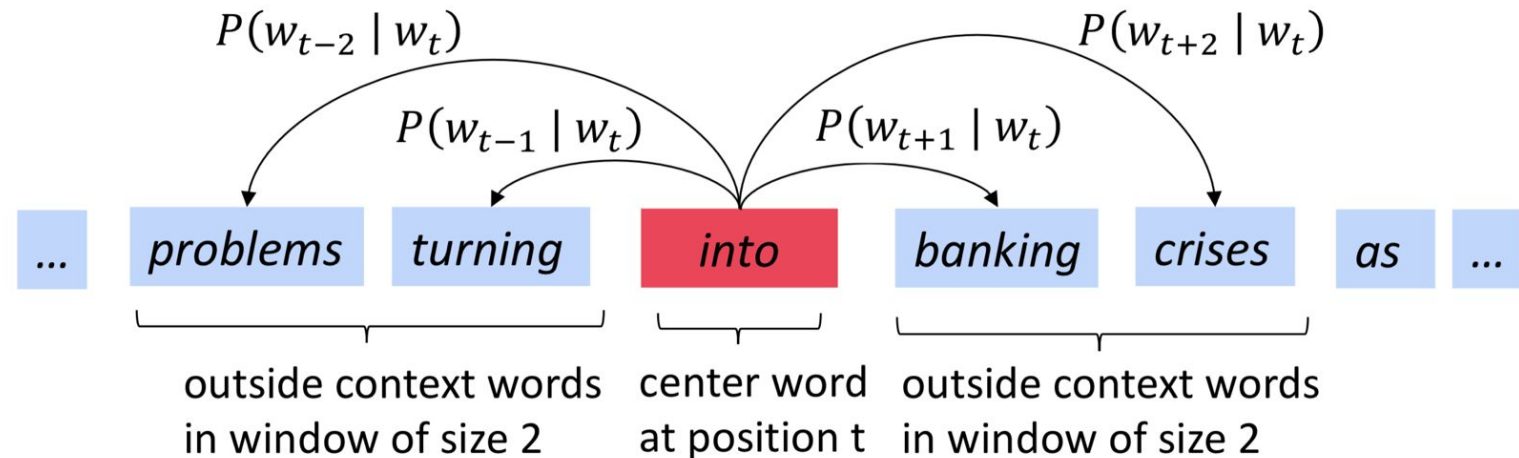
$$\begin{aligned} \text{vec}(d) &= \sum_{w \in V} \text{TF-IDF}(w, d, D) \cdot \text{onehot}(w) \\ &= [\text{TF-IDF}(w_1, d, D), \dots, \text{TF-IDF}(w_n, d, D)] \end{aligned}$$

Что нужно на самом деле?

- Плотные вектора для каждого слова.
- Семантически похожие слова должны иметь похожие вектора.
- Для получения семантически близких векторных представлений можно использовать большие наборы текстов (например, статьи Википедии)
- Как оценить схожесть двух слов при наличии корпуса текстов?
- Похожие слова чаще встречаются в похожих контекстах.
- “You shall know a word by the company it keeps” (c) John Rupert Firth

Алгоритм word2vec

- Выбрать большой набор текстов.
- Определить и зафиксировать словарь слова.
- Случайно инициализировать вектора слов из словаря.
- Для каждой позиции t в тексте с центральным словом s и контекстом o посчитать вероятность встретить слово s в контексте o (или наоборот).
- Обновить вектора для максимизации данной вероятности.
- Повторять, пока вектора не сойдутся.



Как вычислить вероятность при заданных векторах слов?

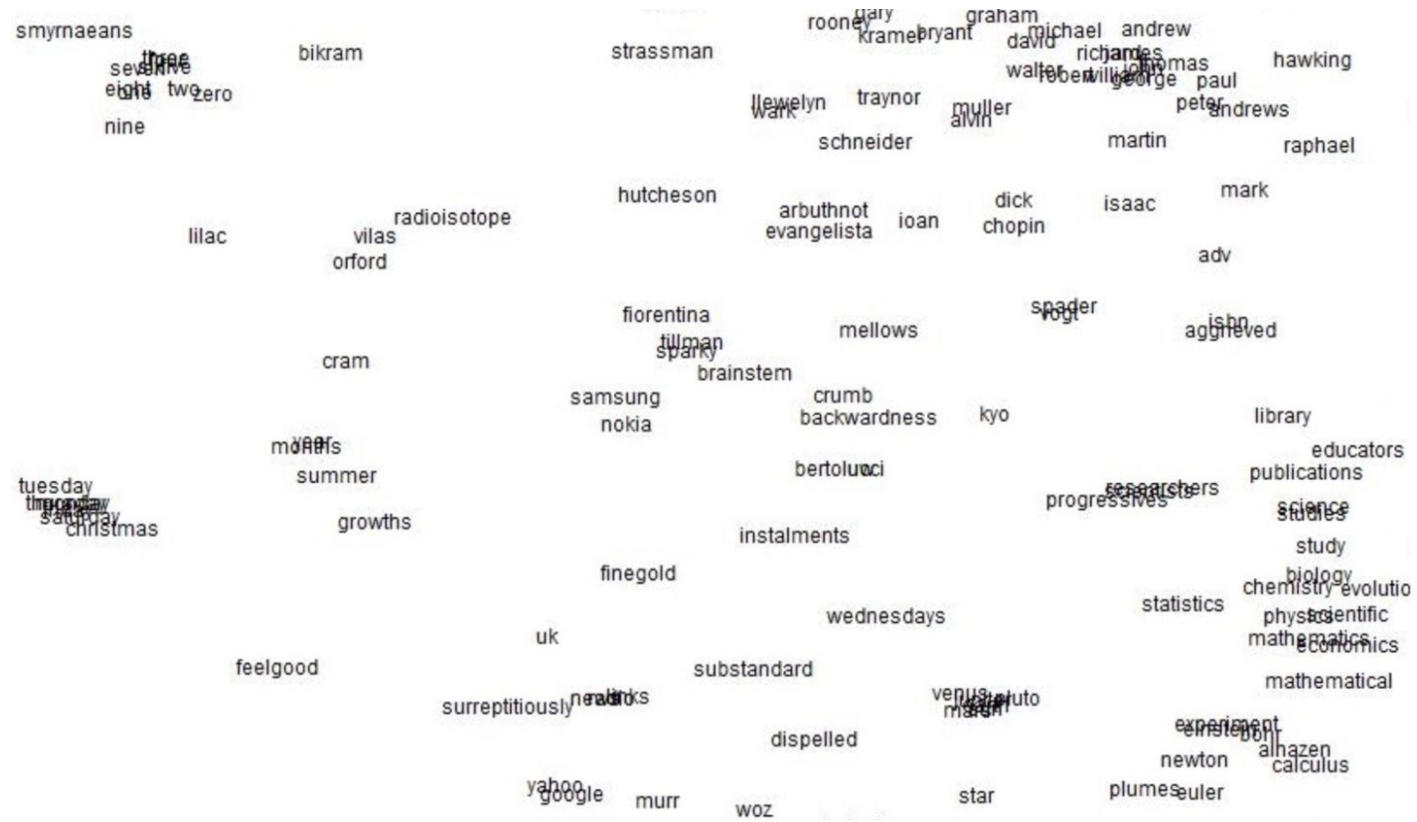
$$P(w_c | w_o) = \frac{e^{w_c \cdot w_o}}{\sum_{w_i \in vocab} e^{w_i \cdot w_o}}$$

$$w_o = \sum_{w_j \in context} w_j$$

$$L = - \sum_{w_i, w_o \in data} \log P(w_c | w_o)$$

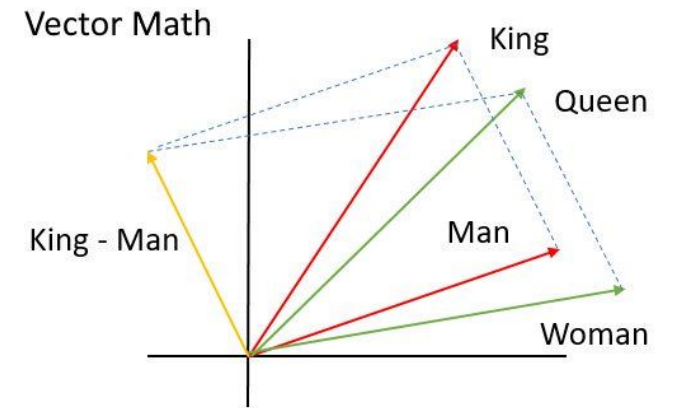
- w_c - вектор центрального слова c
- w_o - вектор контекста o
- $vocab$ - выбранный словарь слов
- $context$ - central word window

Свойства word2vec



Свойства word2vec

- King - Man + Woman \approx Queen
- London - England + Russia \approx Moscow (general domain)
- London - England + Russia \approx Kremlin (news)
- Surgeon - Man + Woman \approx Nurse
- Cosmetics - Woman + Man \approx Pharmaceuticals
- Feminism - Woman + Man \approx Conservatism



GloVe

- word2vec использует только локальную информацию (контекст в несколько слов).
- word2vec не учитывает информацию об общей взаимной встречаемости слов.
- Алгоритм можно улучшить, добавив глобальную статистику по всем документам

$$J = \sum_{i,j=1}^V f(X_{ij}) (w_i^T \tilde{w}_j + b_i + \tilde{b}_j - \log X_{ij})^2$$

- X_{ij} - количество появлений слова j в контексте i ,
- b_i и \tilde{b}_j - коэффициенты смещения для симметризации,
- f - взвешивающая функция для балансирования влияния часто и редко встречающихся слов.

Постановка задачи векторизации текста

Ближайшие соседи слова “frog”:

- Frogs
- Toad
- Litoria
- Leptodactylidae
- Rana
- Lizard
- Eleutherodactylus



3. litoria



4. leptodactylidae



5. rana



7. eleutherodactylus



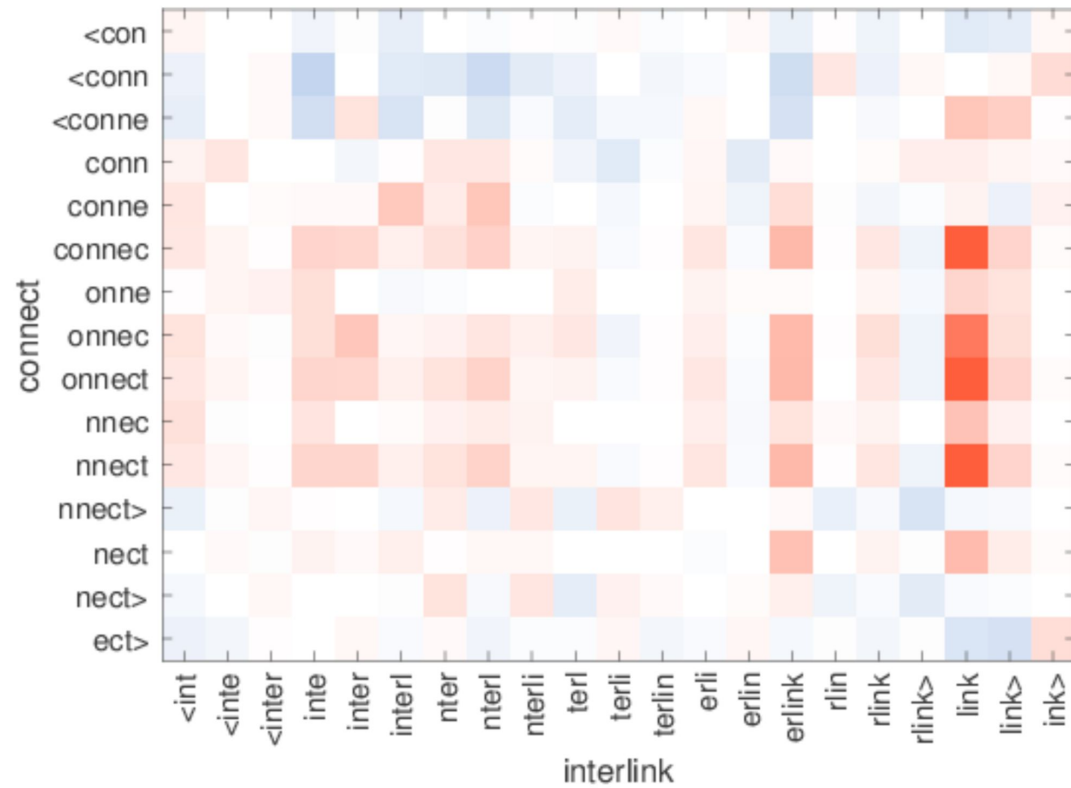
Проблема с out-of-vocabulary словами

- Как выбрать словарь V ?
- Как поступать с out-of-vocabulary словами в тексте?
- Какими именно векторами обозначать out-of-vocabulary слова?

Алгоритм fastText

- Вместо использования словаря слов, используем словарь n-grams.
- n-gram - сочетание n символов.
- Инициализировать векторные представления для n-grams.
- Получить векторное представление слова, сложив вектора входящих в него n-grams.
- Обучить векторные представления методом схожим с word2vec
- fastText n-gram ($n = 3$) представление слова “where”:
 - $\langle wh, whe, her, ere, re \rangle, \langle where \rangle$

Матрица схожести n-grams ($n=6$) алгоритма fastText





Классификация ТЕКСТОВ

Постановка задачи классификации текстов

Пусть X — множество текстов, Y — множество классов.

Существует неизвестная целевая зависимость - отображение $y^*: X \rightarrow Y$, значения которой известны только на объектах конечной обучающей выборки $X^m = \{ (x_1, y_1), \dots, (x_m, y_m) \}$.

Требуется построить алгоритм $a: X \rightarrow Y$, способный классифицировать произвольный объект x из X .

Основные методы решения задачи классификации текстов

- Подход “на правилах” (rule-based or dictionary-based)
- Алгоритмы машинного обучения (Machine Learning)
 - в основном, на векторных представлениях текстов целиком
- Нейронные модели глубокого обучения (Deep Learning)
 - char-level - векторные представления символов
 - *token-level* - векторные представления токенов (слов, пунктуации)
 - sentence-level - векторные представления текстов целиком

Свёрточные нейронные сети

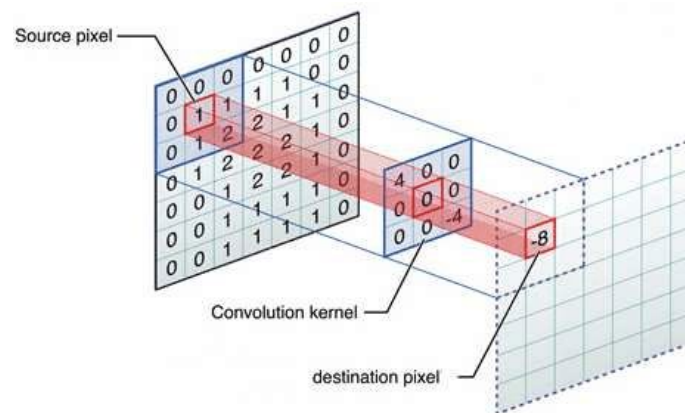
- Свёрточный слой (convolutional layer) включает в себя определенное количество фильтров, ядра (kernel) свёртки которых преобразуют предыдущий слой по фрагментам, суммируя результаты матричного произведения для каждого фрагмента.
- Слой субдискретизация (pooling layer) необходим для нелинейного уплотнения признаков, для избавления от излишне подробных признаков.

1	1	1	0	0
0	1	1	1	0
0	0	1	1	1
0	0	1	1	0
0	1	1	0	0

Image

4		

Convolved Feature



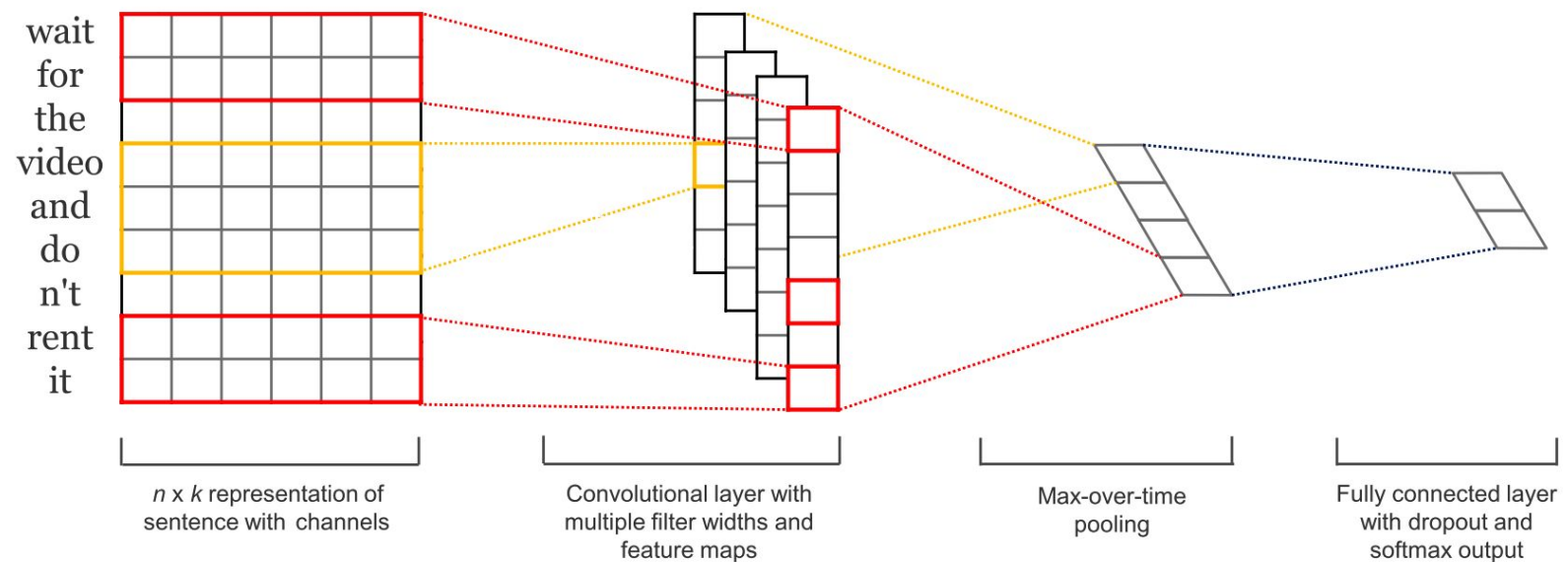
1	4	5	3
2	3	4	7
4	6	2	3
6	3	1	1

4	7
6	3

Max pooling with filter 2x2 and stride of size 2

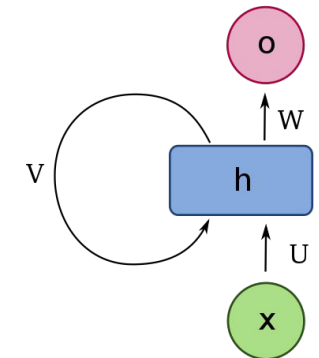
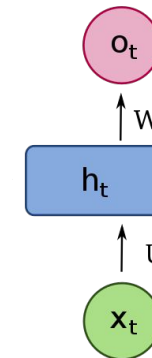
Свёрточные нейронные сети для текстов

- 2D представление текста - потоковое векторное представление размера $n \times k$.
- 1D convolution с f фильтрами возвращает $n \times f$ представление.
- Over-time pooling для ухода от количества токенов.



Рекуррентные нейронные сети

- Текст - это последовательность токенов (временных шагов).
- Для работы с последовательностями необходимо использовать память.
- Переиспользование скрытого состояния с предыдущего временного шага (токена) может решить эту проблему.



Заключение

- word2vec, GloVe, fastText для векторизации текста
- word2vec, GloVe с фиксированным словарем
- fastText работает на уровне n-gram

- Свёрточные нейронные сети для классификации коротких текстов
- Рекуррентные нейронные сети зачастую достигают более высокого качества

DeepPavlov
docs.deeppavlov.ai

Business solutions, support & partnerships
iPavlov.ai

Исследования и разработки выполнены при поддержке Фонда поддержки проектов Национальной технологической инициативы и ПАО "Сбербанк". Идентификатор проекта 0000000007417F630002.

