

Word Embeddings

Sergey Aksenov

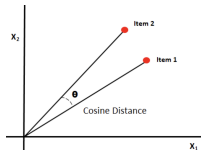
Higher School of Economics

19 сентября 2022 г.

Метрики

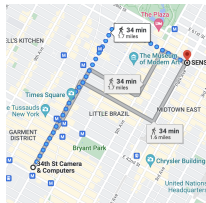
Cosine similarity

$$\text{Sim}(p, q) = \frac{\sum_i p_i q_i}{\sqrt{\sum_i p_i^2} \sqrt{\sum_i q_i^2}}$$



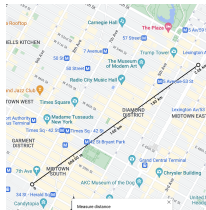
Manhattan distance

$$d(p, q) = \sum_i |p_i - q_i|$$



Euclidean distance

$$d(p, q) = \sqrt{\sum_i (p_i - q_i)^2}$$



One-hot кодирование

Абрикос Апельсин

1
0
0
0

W_1

0
1
0
0

W_2

...

Яблоко

0
0
0
1

$W_{|V|}$

Преимущества

- ▶ Простой способ получить какие-то вектора
- ▶
- ▶

Недостатки

- ▶ Длина векторов не фиксирована
- ▶ Векторы ортогональны
- ▶ Не сохраняет смысл слов

Сингулярное разложение

- ▶ $d \in D, |D| = m$ — документы
- ▶ $w \in V, |V| = n$ — слова
- ▶ $M_{ij} = \#w_i \in d_j$

The diagram illustrates the Singular Value Decomposition (SVD) of matrix M . It shows the equation $M = U \Sigma V^T$ using matrix dimensions and symbols.

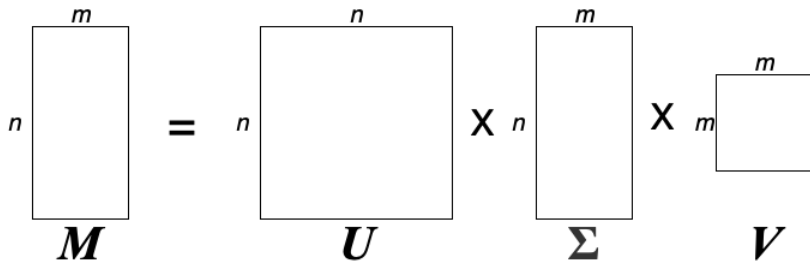
- Matrix M is a vertical rectangle with height n and width m .
- Matrix U is a square with side length n .
- Matrix Σ is a vertical rectangle with height n and width m .
- Matrix V is a square with side length m .

The matrices are arranged in the sequence $M = U \Sigma V^T$, with the dimensions of each matrix indicated by labels n and m next to the corresponding rectangles.

Сингулярное разложение

Сингулярное разложение матрицы M

$$M = U\Sigma V^T$$



Сингулярное разложение

- ▶ Снижение размерности: выберем k наибольших сингулярных чисел в Σ
- ▶ Выберем k первых столбцов и строк в U и $V^T - U_k$ и V_k^T

$$M = U \Sigma V$$

Сингулярное разложение

Что стало лучше

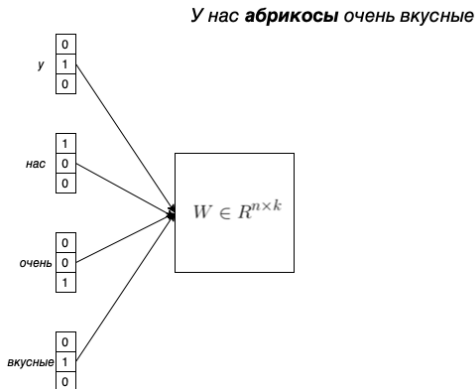
- ▶ Векторы имеют фиксированную длину
- ▶ Векторы не ортогональны друг другу
- ▶ Семантическая близость как-то учитывается

Недостатки

- ▶ Нужно в явном виде работать с большой разреженной матрицей
- ▶ Добавление новых слов или документов невозможно без построения нового разложения

Модель CBOW

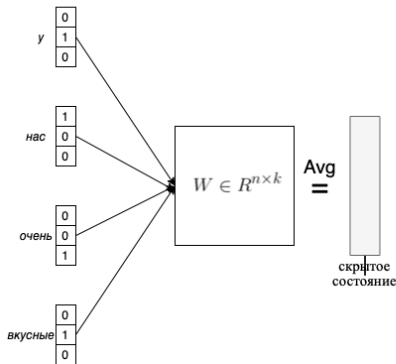
- ▶ Разобьем корпус текста на контексты размера $2C + 1$
- ▶ Вход нейронной сети: 2 вектора контекста, размерность каждого вектра - n



Модель CBOW

- ▶ Применим ко входу преобразование W и усредним результаты
- ▶ Получим k -мерный вектор скрытого состояния
$$h = \frac{1}{4} W^T (a_{i-2} + a_{i-1} + a_{i+1} + a_{i+2})$$

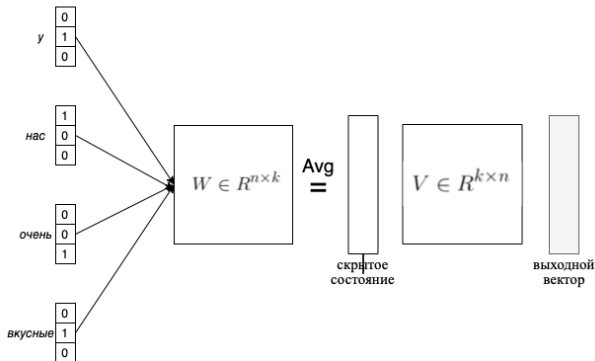
У нас **абрикосы** очень вкусные



Модель CBOW

- ▶ Умножим скрытое состояние на вторую матрицу и получим n -мерный выходной вектор $b = V^T h$

У нас **абрикосы** очень вкусные



Модель CBOW

Оптимизационная задача

- ▶ Для i -го окна по контексту c_i предсказать центральное слово w_i :

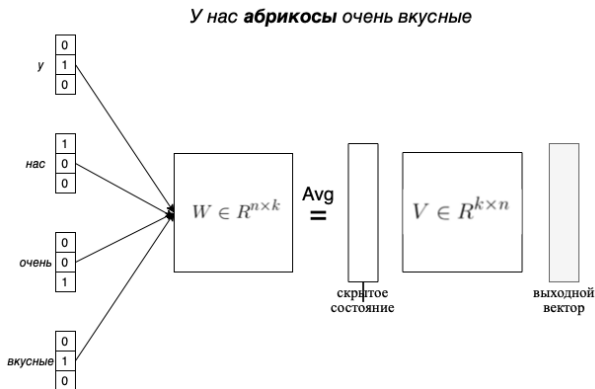
$$\sum_i \log p(w_i | c_i) \rightarrow \max_{W, V}$$

- ▶ Выход последнего слоя сети: $\text{softmax}(b) = (z_1, \dots, z_j)$, где

$$z_j = p(w_j | c_i) = \frac{e^{b_j}}{\sum_k e^{b_k}}$$

Модель CBOW

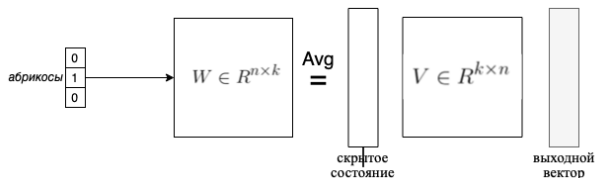
- ▶ Получаем две матрицы весов: W и V
- ▶ Обычно матрицу W считают матрицей искомых векторов



Модель Skip-gram

- ▶ Разобьем корпус на контексты размера $2C + 1$
- ▶ Вход нейронной сети: n –мерный one-hot вектор центрального слова
- ▶ Вектор скрытого состояния $h = W^T x$

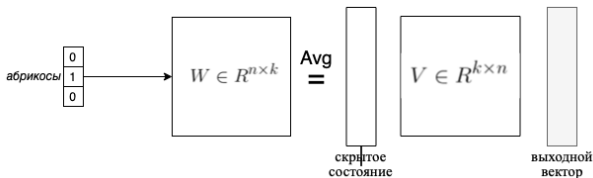
У нас **абрикосы** очень вкусные



Модель Skip-gram

- ▶ Выход нейронной сети: вероятностное распределение слова контекста при условии центрального слова
- ▶ Выходной вектор $b = V^T h$

У нас **абрикосы** очень вкусные



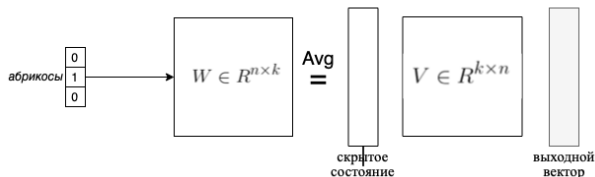
Модель Skip-gram

Оптимизационная задача

- ▶ Для i -го центрального слова предсказать слово контекста c_j

$$\sum_i \sum_j \log(c_j | w_i) \rightarrow \max_{W, V}$$

У нас **абрикосы** очень вкусные



Обучение моделей Word2Vec

Проблема: вычисление *softmax* на выходном слое требует прохода по всему словарю - $O(n)$ операций, вычислительно трудная задача

Методы сокращения перебора

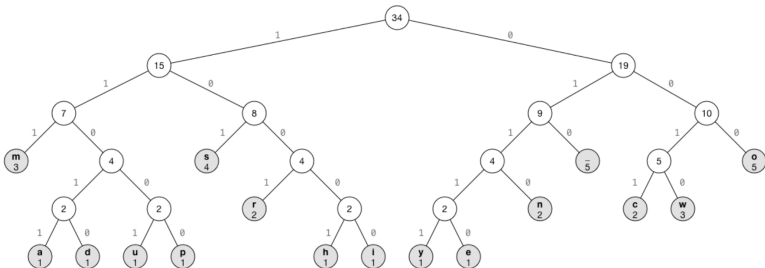
- ▶ Иерархический *softmax*
- ▶ Noise Constructive Estimation (NCE)
- ▶ Negative sampling (NS)

Иерархический softmax

- Будем оценивать только значения, стоящие в позициях предсказываемых слов контекста

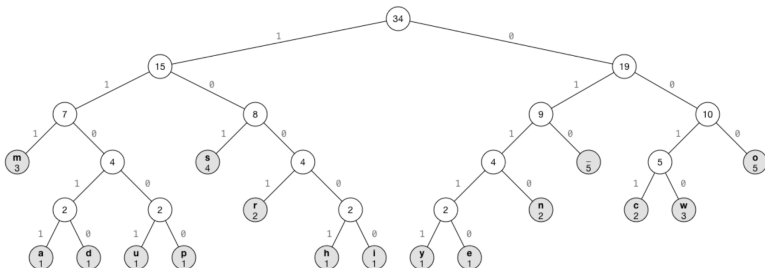
$$\sum_i \sum_j \log(c_j | w_i) \rightarrow \max_{W, V}$$

- Заменяем второй слой сети на дерево Хаффмана



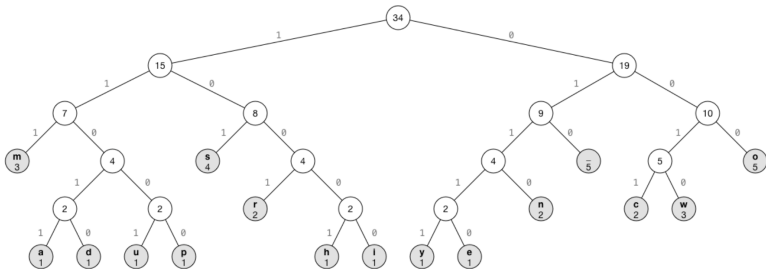
Иерархический softmax

- ▶ Пусть мы хотим оценить вероятность слова "лось"
- ▶ Вероятность перехода из текущей вершины в поддеревья:
 - ▶ $p_0 = \sigma(u_{n(w,0)}^T, h)$, если спускаемся по левой ветке
 - ▶ $p_0 = \sigma(-u_{n(w,0)}^T, h)$, если спускаемся по правой ветке



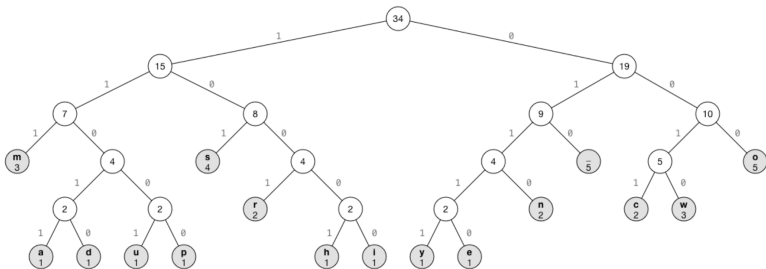
Иерархический softmax

- ▶ В результате мы окажемся в целевом листе
- ▶ Каждый i -й спуск происходил с вероятностью p_i
- ▶ Итоговая вероятность равна $\prod_i p_i$



Иерархический softmax

- ▶ Пройдем по дереву $2C$ раз, получим все вероятности, необходимые для подсчета ошибки и настройки весов
- ▶ Сложность спуска по дереву $O(n \log n)$



Negative sampling

- ▶ Сформулируем задачу skip-gram как задачу бинарной классификации
- ▶ Рассмотрим пару слов (w, s) :
 - ▶ Класс 0: s не входит в контекст w
 - ▶ Класс 1: s входит в контекст w
- ▶ Вероятность класса $p(1|(w, s)) = \frac{1}{1+e^{-\nu_w^t \nu_s}} = \sigma(\nu_w^t, \nu_s)$

Negative sampling

- ▶ Пусть D_1 — множество наблюдаемых пар (w, s)
- ▶ D_0 — множество ненаблюдаемых пар (w, s) , на каждой итерации обучения выбираем случайные пары
- ▶ Функционал правдоподобия

$$L = \sum_{w, s \in D_1} \log(\sigma(\nu_w^t \nu_s)) + \sum_{w, s \in D_0} \log(\sigma(-\nu_w^t \nu_s))$$

Модели word2vec

Обучение

- ▶ Ускорение обучения за счет иерархического softmax и negative sampling
- ▶ Вариант по умолчанию: SGNS (skip-gram with negative sampling)

Модели word2vec

Обучение

- ▶ Ускорение обучения за счет иерархического softmax и negative sampling
- ▶ Вариант по умолчанию: SGNS (skip-gram with negative sampling)

Недостатки

- ▶ Не учитывается морфология
- ▶ Нет возможности работать с OOV-словами

Модель FastText

- ▶ Разобьем слова на последовательности символов — символьные N-граммы
- ▶ Обучим вектора для каждой N-граммы с помощью CBOW или Skip-gram
- ▶ Вектор слова получим усреднением N-грамм