

# Deep Learning. Анализ текстов

## Урок 9

---

Егор Конягин

25 июля 2019 г.

МФТИ & АО "ЦОСиВТ"

1. Повторение
2. Векторизация текста (word embeddings)
3. Как получить адекватный word embedding?

# Повторение

---

В модели RNN ответ, полученный на  $t$ -ом элементе  $x$ , передается далее

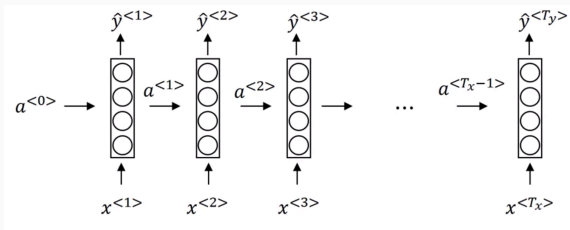


Рис. 1: Простейшая модель RNN. Источник: Andrew Ng

$$a^{<t>} = \sigma_1(w_{aa} \cdot a^{<t-1>} + w_{ax} \cdot x^{<t>} + b_a) \quad a^{<0>} = \vec{0} \quad (1)$$

$$y^{<t>} = \sigma_1(w_{ya} \cdot a^{<t-1>} + b_a) \quad (2)$$

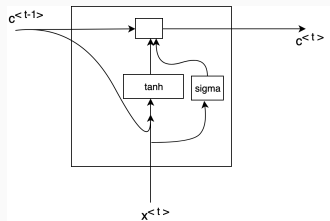


Рис. 2: GRU-ячейка

Обновление происходит только, если  $\Gamma$  примет соответствующее значение (затвор).

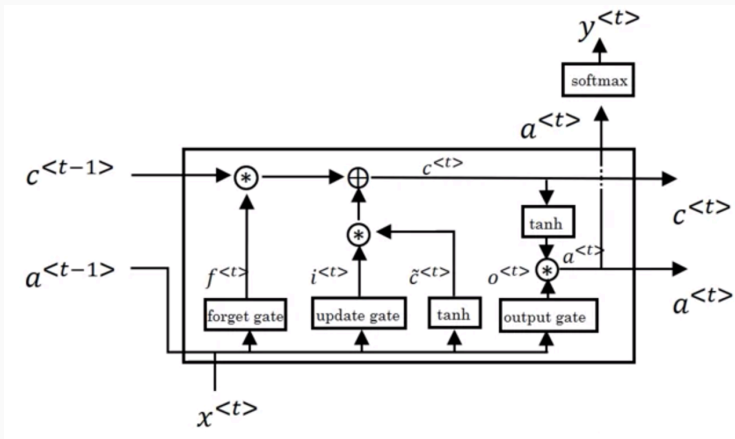


Рис. 3: LSTM-ячейка. Источник: Andrew Ng's classes

# Векторизация текста (word embeddings)

---

Мы уже рассматривали one-hot векторизацию текста и отметили, что такая модель не является оптимальной. **Почему?**

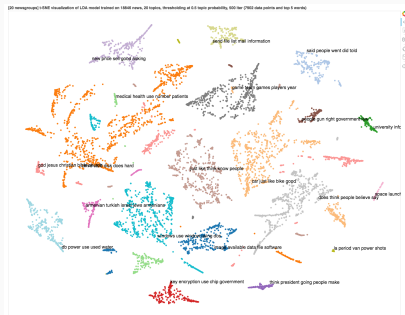
1. векторы слов, по смыслу похожих, ничего общего не имеют;
2. разреженное представление увеличивает затратность ресурсов;

**Решение:** использовать координаты в векторном представлении как значения неких заданных признаков.

Построив такую модель языка, можно использовать обычное скалярное произведение для оценки похожести слов/текста...



Алгоритм t - SNE - это алгоритм понижения размерности для визуализации данных (как и PCA). Однако он, в отличие от PCA, является нелинейным алгоритмом и лучше позволяет визуализировать данные.



**Рис. 4:** Использование t-sne для отображения тем новостей. Источник: SHUAI'S AI & DATA BLOG

Как получить векторное представление и дальше работать с ним?

1. Получить представления, обучаясь на больших объемах текста (миллиард слов).
2. Использовать эти представления для своей задачи с помощью дообучения на маленьком датасете (тысячи слов).

# Word embeddings. Пример

Векторные представления способны понимать признаки, присущие слову. Рассмотрим пример:

Мужчина  $\rightarrow$  женщина = король  $\rightarrow$  кто?

Решение:

$$\vec{e}_{man} - \vec{e}_{woman} = \vec{r}; \quad (3)$$

$$\vec{e}_{king} - \vec{e}_i = \vec{r}_i; \quad (4)$$

Наша задача найти такой  $i'$ , чтобы  $\mathcal{D}(\vec{r}, \vec{r}_{i'}) \rightarrow \min$ . Как выбрать функцию  $\mathcal{D}$ ?

Обычно в анализе текста используют cosine similarity:

$$\mathcal{D}(a, b) = \frac{a^T b}{|a| \cdot |b|}. \quad (5)$$

## Как получить адекватный word embedding?

Положим, мы составляем представления для 10 000 различных слов, а представления - это 300-мерные векторы. Тогда можно составить матрицу для перехода из one-hot представления к более адекватному (низкой размерности).

Эта матрица носит название embedding matrix.

Как получить адекватный word  
embedding?

---

Построение вероятностной модели языка:

1. задать матрицу  $E$  (10000,300) случайным образом;
2. подать выход этой матрицы на вход полносвязной нейросети с softmax-выходом;
3. на вход алгоритму подавать  $m$  слов предложения. Задача алгоритма - предсказать следующее слово (то есть решать задачу многоклассовой классификации);
4. запустить обучение;
5. ЖДЕМ;
6. Готово!

Как и в предыдущем случае, ставится задача обучения с учителем:  $x$  - это т.н. контекст,  $y$  - цель (одно слово).

В данном подходе контекст строится существенно проще:

1. выбирается случайное слово в предложении - это и будет контекстом;
2. в пределах окна фиксированной ширины ( $\pm 5$  слов) выбирается случайное слово - это и будет целевое слово.

Недостатки:

- при использовании словаря большего размера softmax-слой становится огромным;
- случайно выбирая слово для контекста, мы наткнемся на наиболее используемые слова в языке (артикли, предлоги и тд).

Решение: использование иерархического классификатора, а затем softmax-слоя.

# Negative sampling

Идея: сформулировать задачу как бинарную классификацию:

$$(e_i, e_j) \rightarrow \text{classifier} \rightarrow \mathcal{P}(e_i, e_j) \quad (6)$$

Используя данную модель, нам надо также в обучающую выборку добавить заведомо несочетающиеся слова (отсюда название negative sampling).

**Как выбирать слова для сэмплинга?** Ответ: используя частотный анализ (с небольшими модификациями).

$$\mathcal{P}(w_i) = \frac{f(w_i)}{\sum_j f(w_j)} \quad (7)$$

$$\mathcal{P}^*(w_i) = \frac{f(w_i)^{3/4}}{\sum_j f(w_j)^{3/4}} \quad (8)$$

**Важно!** На одно положительное словосочетание должно приходиться к негативных.



Снова рассматриваем контекст  $C$  и целевое слово  $T$ .

В явном виде запишем величину  $X_{ij}$  - кол-во появлений  $i$ -ого слова в контексте  $j$ -ого слова.

В общем случае  $X_{ij} \neq X_{ji}$ .

Постановка задачи:

$$\sum_i \sum_j f(x_{ij})(\theta_i^T \cdot e_j - \log x_{ij})^2 \rightarrow \min_{\theta, e} \quad (9)$$

# Summary