



# Deep Learning School

Физтех-Школа Прикладной математики и информатики (ФПМИ) МФТИ

---

*Some parts of the notebook are almost the copy of [mmta-team course](#). Special thanks to mmta-team for making them publicly available. [Original notebook](#).*

Прочитайте семинар, пожалуйста, для успешного выполнения домашнего задания. В конце ноутка напишите свой вывод. Работа без вывода оценивается ниже.

## ▼ Задача поиска схожих по смыслу предложений

Мы будем ранжировать вопросы [StackOverflow](#) на основе семантического векторного представления

До этого в курсе не было речи про задачу ранжирования, поэтому введем математическую формулировку

## ▼ Задача ранжирования(Learning to Rank)



- $X$  - множество объектов
- $X^l = \{x_1, x_2, \dots, x_l\}$  - обучающая выборка

На обучающей выборке задан порядок между некоторыми элементами, то есть нам известно, что некий объект выборки более релевантный для нас, чем другой:

- $i < j$  - порядок пары индексов объектов на выборке  $X^l$  с индексами  $i$  и  $j$

▼ Задача:

построить ранжирующую функцию  $a : X \rightarrow R$  такую, что

$$i < j \Rightarrow a(x_i) < a(x_j)$$



▼ Embeddings

Будем использовать предобученные векторные представления слов на постах Stack Overflow.  
[A word2vec model trained on Stack Overflow posts](#)

```
#!curl https://zenodo.org/record/1199620/files/S0_vectors_200.bin?download=1
```

```
from gensim.models.keyedvectors import KeyedVectors  
  
wv_embeddings = KeyedVectors.load_word2vec_format("S0_vectors_200.bin", binary=True)
```

## Как пользоваться этими векторами?

Посмотрим на примере одного слова, что из себя представляет embedding

```
word = "dog"  
if word in wv_embeddings:  
    print(wv_embeddings[word].dtype, wv_embeddings[word].shape)  
  
float32 (200,)
```

```
print(f"Num of words: {len(wv_embeddings.index_to_key)}")
```

```
Num of words: 1787145
```

Найдем наиболее близкие слова к слову `dog`:

## Вопрос 1:

- Входит ли слово `cat` в топ-5 близких слов к слову `dog`? Какое место оно занимает?

```
def most_similar(word: str, n: int):  
    return wv_embeddings.most_similar(word)[:n]
```

```
most_similar("dog", 5)
```

```
[('animal', 0.8564179539680481),  
 ('dogs', 0.7880867123603821),  
 ('mammal', 0.7623804211616516),  
 ('cats', 0.7621253728866577),  
 ('animals', 0.7607938647270203)]
```

**Ваш ответ:** 'cat' не входит в топ-5 близайших слов к 'dog'.

## Векторные представления текста

Перейдем от векторных представлений отдельных слов к векторным представлениям вопросов, как к **среднему** векторов всех слов в вопросе. Если для какого-то слова нет предобученного вектора, то его нужно пропустить. Если вопрос не содержит ни одного известного слова, то нужно вернуть нулевой вектор.

```
import numpy as np
import re

# you can use your tokenizer
# for example, from nltk.tokenize import WordPunctTokenizer
class MyTokenizer:
    def __init__(self):
        pass

    def tokenize(self, text):
        return re.findall("\w+", text)

tokenizer = MyTokenizer()
```

```
import nltk
from nltk.tokenize import TreebankWordTokenizer
from nltk.stem import WordNetLemmatizer
from nltk.corpus import wordnet

nltk.download("averaged_perceptron_tagger_eng")
nltk.download("wordnet")
nltk.download("omw-1.4")

class NLTKTokenizer:
    def __init__(self):
        self.tokenizer = TreebankWordTokenizer()
        self.lemmatizer = WordNetLemmatizer()

    def get_wordnet_pos(self, treebank_tag):
        if treebank_tag.startswith("J"):
            return wordnet.ADJ
        elif treebank_tag.startswith("V"):
            return wordnet.VERB
        elif treebank_tag.startswith("N"):
            return wordnet.NOUN
        elif treebank_tag.startswith("R"):
            return wordnet.ADV
        else:
            return wordnet.NOUN

    def tokenize(self, text):
        tokens = self.tokenizer.tokenize(text)
        pos_tags = nltk.pos_tag(tokens)
        lemmas = [
            self.lemmatizer.lemmatize(tok.lower(), self.get_wordnet_pos(pos))
            for tok, pos in pos_tags
        ]
        return lemmas
```

```
nltk_tok = NLTKTokenizer()
```

```
[nltk_data] Downloading package averaged_perceptron_tagger_eng to
[nltk_data]      /Users/arsknz/nltk_data...
[nltk_data] Package averaged_perceptron_tagger_eng is already up-to-
[nltk_data]      date!
[nltk_data] Downloading package wordnet to /Users/arsknz/nltk_data...
[nltk_data] Package wordnet is already up-to-date!
[nltk_data] Downloading package omw-1.4 to /Users/arsknz/nltk_data...
[nltk_data] Package omw-1.4 is already up-to-date!
```

```
tokenizer.tokenize("I love neural networks")
```

```
['I', 'love', 'neural', 'networks']
```

```
def question_to_vec(question, embeddings=wv_embeddings, tokenizer=tokenizer, dim=200
    """
        question: строка
        embeddings: наше векторное представление
        dim: размер любого вектора в нашем представлении

        return: векторное представление для вопроса
    """
    tokens = tokenizer.tokenize(question)

    sent = []

    for i in tokens:
        if i in embeddings:
            sent.append(embeddings[i])
    if len(sent) != 0:
        return np.mean(sent, axis=0)
    else:
        return np.zeros(
            200,
        )
```

```
def cos_sim(nums1, nums2):
    norm1 = np.linalg.norm(nums1)
    norm2 = np.linalg.norm(nums2)
    if norm1 == 0 or norm2 == 0:
        return 0.0
    return np.dot(nums1, nums2) / (norm1 * norm2)
```

Теперь у нас есть метод для создания векторного представления любого предложения.

## ▼ Вопрос 2:

- Какая третья (с индексом 2) компонента вектора предложения I love neural networks (округлите до 2 знаков после запятой)?

```
# Предложение
question = "I love neural networks"

question_to_vec(question)[2]

-1.2854122
```

-1.2854122

## ▼ Оценка близости текстов

Представим, что мы используем идеальные векторные представления слов. Тогда косинусное расстояние между дублирующими предложениями должно быть меньше, чем между случайно взятыми предложениями.

Сгенерируем для каждого из  $N$  вопросов  $R$  случайных отрицательных примеров и примешаем к ним также настоящие дубликаты. Для каждого вопроса будем ранжировать с помощью нашей модели  $R + 1$  примеров и смотреть на позицию дубликата. Мы хотим, чтобы дубликат был первым в ранжированном списке.

### Hits@K

Первой простой метрикой будет количество корректных попаданий для какого-то  $K$ :

$$\text{Hits@K} = \frac{1}{N} \sum_{i=1}^N [rank\_q'_i \leq K],$$

$$[x < 0] \equiv \begin{cases} 1, & x < 0 \\ 0, & x \geq 0 \end{cases}$$

- индикаторная функция
- $q_i$  -  $i$ -ый вопрос
- $q'_i$  - его дубликат
- $rank\_q'_i$  - позиция дубликата в ранжированном списке ближайших предложений для вопроса  $q_i$ .

Hits@K измеряет долю вопросов, для которых правильный ответ попал в топ- $K$  позиций среди отранжированных кандидатов.

### DCG@K

Второй метрикой будет упрощенная DCG метрика, учитывающая порядок элементов в списке путем домножения релевантности элемента на вес равный обратному логарифму номера позиции::

$$\text{DCG@K} = \frac{1}{N} \sum_{i=1}^N \frac{1}{\log_2(1 + rank\_q'_i)} \cdot [rank\_q'_i \leq K],$$

С такой метрикой модель штрафуется за большой ранк корректного ответа.

DCG@K измеряет качество ранжирования, учитывая не только факт наличия правильного ответа в топ-K, но и *его точную позицию*.



▼ Пример оценок

Вычислим описанные выше метрики для игрушечного примера. Пусть

- $N = 1, R = 3$
- "Что такое python?" - вопрос  $q_1$
- "Что такое язык python?" - его дубликат  $q'_i$

Пусть модель выдала следующий ранжированный список кандидатов:

1. "Как изучить c++?"
2. "Что такое язык python?"
3. "Хочу учить Java"
4. "Не понимаю Tensorflow"

$$\Rightarrow \text{rank}_{q'_i} = 2$$

Вычислим метрику  $\text{Hits}@K$  для  $K = 1, 4$ :

- [ $K = 1$ ]  $\text{Hits}@1 = [\text{rank}_{q'_i} \leq 1]$

Проверяем условие  $\text{rank}_{q'_i} \leq 1$ : **условие неверно.**

Следовательно,  $[\text{rank}_{q'_i} \leq 1] = 0$ .

- [ $K = 4$ ]  $\text{Hits}@4 = [\text{rank}_{q'_i} \leq 4] = 1$

Проверяем условие  $\text{rank}_{q'_i} \leq 4$ : **условие верно.**

Вычислим метрику  $\text{DCG}@K$  для  $K = 1, 4$ :

- [ $K = 1$ ]  $\text{DCG}@1 = \frac{1}{\log_2(1+2)} \cdot [2 \leq 1] = 0$
- [ $K = 4$ ]  $\text{DCG}@4 = \frac{1}{\log_2(1+2)} \cdot [2 \leq 4] = \frac{1}{\log_2 3}$

▼ **Вопрос 3:**

- Вычислите DCG@10, если  $rank_{q'_i} = 9$  (округлите до одного знака после запятой)

1 / np.log2(1 + 9) \* (9 <= 10)

0.3010299956639812

▼ ~ 0.3

▼ Более сложный пример оценок

Рассмотрим пример с  $N > 1$ , где  $N = 3$  (три вопроса) и для каждого вопроса заданы позиции их дубликатов. Вычислим метрики Hits@K для разных значений  $K$ .

- $N = 3$ : Три вопроса ( $q_1, q_2, q_3$ ).
- Для каждого вопроса известна позиция его дубликата ( $rank_{q'_i}$ ):
  - $rank_{q'_1} = 2$ ,
  - $rank_{q'_2} = 5$ ,
  - $rank_{q'_3} = 1$ .

Мы будем вычислять Hits@K для  $K = 1, 5$ .

**Для  $K = 1$ :**

Подставим значения:

$$\text{Hits}@1 = \frac{1}{3} \cdot ([\text{rank}_{q'_1} \leq 1] + [\text{rank}_{q'_2} \leq 1] + [\text{rank}_{q'_3} \leq 1]).$$

Проверяем условие  $\text{rank}_{q'_i} \leq 1$  для каждого вопроса:

- $\text{rank}_{q'_1} = 2 \rightarrow 2 \not\leq 1 \rightarrow 0$ ,
- $\text{rank}_{q'_2} = 5 \rightarrow 5 \not\leq 1 \rightarrow 0$ ,
- $\text{rank}_{q'_3} = 1 \rightarrow 1 \leq 1 \rightarrow 1$ .

Сумма:

$$\text{Hits}@1 = \frac{1}{3} \cdot (0 + 0 + 1) = \frac{1}{3}.$$

$$\boxed{\text{Hits}@1 = \frac{1}{3}}.$$

**Для  $K = 5$ :**

Подставим значения:

$$\text{Hits}@5 = \frac{1}{3} \cdot ([\text{rank}_{q'_1} \leq 5] + [\text{rank}_{q'_2} \leq 5] + [\text{rank}_{q'_3} \leq 5]).$$

Проверяем условие  $\text{rank}_{q'_i} \leq 5$  для каждого вопроса:

- $\text{rank}_{q'_1} = 2 \rightarrow 2 \leq 5 \rightarrow 1$ ,
- $\text{rank}_{q'_2} = 5 \rightarrow 5 \leq 5 \rightarrow 1$ ,
- $\text{rank}_{q'_3} = 1 \rightarrow 1 \leq 5 \rightarrow 1$ .

Сумма:

$$\text{Hits@5} = \frac{1}{3} \cdot (1 + 1 + 1) = 1.$$

$$\boxed{\text{Hits@5} = 1}.$$

Теперь вычислим метрику **DCG@K** для того же примера, где  $N = 3$  (три вопроса), и для каждого вопроса известна позиция его дубликата ( $\text{rank}_{q'_i}$ ):

- $\text{rank}_{q'_1} = 2$ ,
- $\text{rank}_{q'_2} = 5$ ,
- $\text{rank}_{q'_3} = 1$ .

Мы будем вычислять **DCG@K** для  $K = 1, 5$ .

---

**Для  $K = 1$ :** Подставим значения:

$$\text{DCG}@1 = \frac{1}{3} \cdot \left( \frac{1}{\log_2(1 + \text{rank}_{q'_1})} \cdot [\text{rank}_{q'_1} \leq 1] + \frac{1}{\log_2(1 + \text{rank}_{q'_2})} \cdot [\text{rank}_{q'_2} \leq 1] + \frac{1}{\log_2(1 + \text{rank}_{q'_3})} \cdot [\text{rank}_{q'_3} \leq 1] \right)$$

Проверяем условие  $\text{rank}_{q'_i} \leq 1$  для каждого вопроса:

- $\text{rank}_{q'_1} = 2 \rightarrow 2 \not\leq 1 \rightarrow 0$ ,
- $\text{rank}_{q'_2} = 5 \rightarrow 5 \not\leq 1 \rightarrow 0$ ,
- $\text{rank}_{q'_3} = 1 \rightarrow 1 \leq 1 \rightarrow 1$ .

Сумма:

$$\text{DCG}@1 = \frac{1}{3} \cdot (0 + 0 + 1) = \frac{1}{3}.$$

$$\boxed{\text{DCG}@1 = \frac{1}{3}}.$$

**Для  $K = 5$ :** Подставим значения:

$$\text{DCG}@5 = \frac{1}{3} \cdot \left( \frac{1}{\log_2(1 + \text{rank}_{q'_1})} \cdot [\text{rank}_{q'_1} \leq 5] + \frac{1}{\log_2(1 + \text{rank}_{q'_2})} \cdot [\text{rank}_{q'_2} \leq 5] + \frac{1}{\log_2(1 + \text{rank}_{q'_3})} \cdot [\text{rank}_{q'_3} \leq 5] \right)$$

Проверяем условие  $\text{rank}_{q'_i} \leq 5$  для каждого вопроса:

- $\text{rank}_{q'_1} = 2 \rightarrow 2 \leq 5 \rightarrow 1$ ,
- $\text{rank}_{q'_2} = 5 \rightarrow 5 \leq 5 \rightarrow 1$ ,
- $\text{rank}_{q'_3} = 1 \rightarrow 1 \leq 5 \rightarrow 1$ .

Сумма:

$$\text{DCG}@5 = \frac{1}{3} \cdot (0.631 + 0.387 + 1) = \frac{1}{3} \cdot 2.018 \approx 0.673.$$

$\boxed{\text{DCG@5} \approx 0.673}$ .

▼ **Вопрос 4:**

- Найдите максимум  $\boxed{\text{Hits@47} - \text{DCG@1}}$ ?

```
print(
(
    np.mean([[i <= 47] for i in [2, 5, 1]]),
    np.mean([1 / np.log2(1 + i) * (i <= 1) for i in [2, 5, 1]]),
)
)
print(
(
    np.mean([[i <= 47] for i in [2, 5, 1]])
    - np.mean([1 / np.log2(1 + i) * (i <= 1) for i in [2, 5, 1]]),
)
)
(1.0, 0.3333333333333333)
(0.6666666666666667,)
```

▼ HITS\_COUNT и DCG\_SCORE

Каждая функция имеет два аргумента:  $dup\_ranks$  и  $k$ .

$dup\_ranks$  является списком, который содержит рейтинги дубликатов (их позиции в ранжированном списке).

К примеру для "Что такое язык python?"  $dup\_ranks = [2]$ .

```
def hits_count(dup_ranks, k):
    """
    dup_ranks: list индексов дубликатов
    k: пороговое значение для ранга
    result: вернуть Hits@k
    """
    # Подсчитываем количество дубликатов, чей ранг <= k
    hits_value = np.mean([i <= k for i in dup_ranks])
    return hits_value
```

```
dup_ranks = [2]

k = 1
hits_value = hits_count(dup_ranks, k)
print(f"Hits@1 = {hits_value}")

k = 4
```

```
hits_value = hits_count(dup_ranks, k)
print(f"Hits@4 = {hits_value}")

Hits@1 = 0.0
Hits@4 = 1.0
```

```
def dcg_score(dup_ranks, k):
    """
    dup_ranks: list индексов дубликатов
    k: пороговое значение для ранга
    result: вернуть DCG@k
    """
    dcg_value = np.mean([1 / np.log2(1 + i + 1) * (i + 1 <= k) for i in dup_ranks])
    return dcg_value
```

```
# Пример списка позиций дубликатов
dup_ranks = [2]

# Вычисляем DCG@1
dcg_value = dcg_score(dup_ranks, k=1)
print(f"DCG@1 = {dcg_value:.3f}")

# Вычисляем DCG@4
dcg_value = dcg_score(dup_ranks, k=4)
print(f"DCG@10 = {dcg_value:.3f}")

DCG@1 = 0.000
DCG@10 = 0.500
```

Протестируем функции. Пусть  $N = 1$ , то есть один эксперимент. Будем искать копию вопроса и оценивать метрики.

```
import pandas as pd

copy_answers = [
    "How does the catch keyword determine the type of exception that was thrown",
]

# наши кандидаты
candidates_ranking = [
    [
        "How Can I Make These Links Rotate in PHP",
        "How does the catch keyword determine the type of exception that was thrown",
        "NSLog array description not memory address",
        "PECL_HTTP not recognised php ubuntu",
    ],
]

# dup_ranks – позиции наших копий, так как эксперимент один, то этот массив длины 1
dup_ranks = [1]

# вычисляем метрику для разных k
```

```
print("Ваш ответ HIT:", [hits_count(dup_ranks, k) for k in range(1, 5)])
print("Ваш ответ DCG:", [round(dcg_score(dup_ranks, k), 5) for k in range(1, 5)])
```

```
Ваш ответ HIT: [1.0, 1.0, 1.0, 1.0]
Ваш ответ DCG: [0.0, 0.63093, 0.63093, 0.63093]
```

У вас должно получиться

```
# correct_answers – метрика для разных k
correct_answers = pd.DataFrame(
    [[0, 1, 1, 1], [0, 1 / (np.log2(3)), 1 / (np.log2(3)), 1 / (np.log2(3))]],
    index=["HITS", "DCG"],
    columns=range(1, 5),
)
correct_answers
```

	1	2	3	4
<b>HITS</b>	0	1.00000	1.00000	1.00000
<b>DCG</b>	0	0.63093	0.63093	0.63093

## ▼ Данные

[arxiv link](#)

`train.tsv` - выборка для обучения.

В каждой строке через табуляцию записаны: **<вопрос>, <похожий вопрос>**

`validation.tsv` - тестовая выборка.

В каждой строке через табуляцию записаны: **<вопрос>, <похожий вопрос>, <отрицательный пример 1>, <отрицательный пример 2>, ...**

Считайте данные.

```
def read_corpus(filename):
    data = []
    with open(filename, encoding="utf-8") as file:
        for line in file:
            data.append(line.strip().split("\t"))
    return data
```

Нам понадобится только файл validation.

```
validation_data = read_corpus("./data/validation.tsv")
```

Кол-во строк

```
len(validation_data)
```

3760

Размер нескольких первых строк

```
for i in range(25):
    print(i + 1, len(validation_data[i]))

1 1001
2 1001
3 1001
4 1001
5 1001
6 1001
7 1001
8 1001
9 1001
10 1001
11 1001
12 1001
13 1001
14 1001
15 1001
16 1001
17 1001
18 1001
19 1001
20 1001
21 1001
22 1001
23 1001
24 1001
25 1001
```

## ▼ Ранжирование без обучения

Реализуйте функцию ранжирования кандидатов на основе косинусного расстояния. Функция должна по списку кандидатов вернуть отсортированный список пар (позиция в исходном списке кандидатов, кандидат). При этом позиция кандидата в полученном списке является его рейтингом (первый - лучший). Например, если исходный список кандидатов был [a, b, c], и самый похожий на исходный вопрос среди них - с, затем а, и в конце б, то функция должна вернуть список [(2, c), (0, a), (1, b)].

```
def rank_candidates(question, candidates, embeddings, tokenizer, dim=200):
    """
    question: строка
    candidates: массив строк(кандидатов) [a, b, c]
    result: пары (начальная позиция, кандидат) [(2, c), (0, a), (1, b)]
    """
    sim = []
```

```

for pos, data in enumerate(candidates):
    sim.append(
        [
            pos,
            data,
            cos_sim(
                question_to_vec(
                    question, embeddings=embeddings, tokenizer=tokenizer
                ),
                question_to_vec(data, embeddings=embeddings, tokenizer=tokenizer
            )),
        ]
    )
sim = sorted(sim, key=lambda x: x[2])[::-1]
return [i[:2] for i in sim]

```

Протестируйте работу функции на примерах ниже. Пусть  $N = 2$ , то есть два эксперимента

```

questions = ["converting string to list", "Sending array via Ajax fails"]

candidates = [
    [
        "Convert Google results object (pure js) to Python object", # первый эксперимент
        "C# create cookie from string and send it",
        "How to use jQuery AJAX for an outside domain?",
    ],
    [
        "Getting all list items of an unordered list in PHP", # второй эксперимент
        "WPF– How to update the changes in list item of a list",
        "select2 not displaying search results",
    ],
]

```

```

for question, q_candidates in zip(questions, candidates):
    ranks = rank_candidates(question, q_candidates, wv_embeddings, tokenizer)
    print(ranks)
    print()

[[1, 'C# create cookie from string and send it'], [0, 'Convert Google results object'],
 [[1, 'WPF– How to update the changes in list item of a list'], [0, 'Getting all list

```

Для первого эксперимента вы можете полностью сравнить ваши ответы и правильные ответы. Но для второго эксперимента два ответа на кандидаты будут **скрыты**(\*)

```

# должно вывести
"""
results = [[(1, 'C# create cookie from string and send it'),
            (0, 'Convert Google results object (pure js) to Python object'),
            (2, 'How to use jQuery AJAX for an outside domain?')],
```

```
[(*, 'Getting all list items of an unordered list in PHP'), #скрыт
 (*, 'select2 not displaying search results'), #скрыт
 (*, 'WPF- How to update the changes in list item of a list')]] #скрыт
"""

"\nresults = [[(1, 'C# create cookie from string and send it'),\n          (0,
 'Convert Google results object (pure js) to Python object'),\n          (2, 'How
 to use jQuery AJAX for an outside domain?')],\n          [(*, 'Getting all list
 items of an unordered list in PHP'), #скрыт\n          (*, 'select2 not displaying
 search results'), #скрыт\n          (*, 'WPF- How to update the changes in list
 item of a list')]] #скрыт\n          "
```

Последовательность начальных индексов вы должны получить для эксперимента 1 1, 0, 2.

### Вопрос 5:

- Какую последовательность начальных индексов вы получили для эксперимента 2(перечисление без запятой и пробелов, например, 102 для первого эксперимента?

▼ 102

Теперь мы можем оценить качество нашего метода. Запустите следующие два блока кода для получения результата. Обратите внимание, что вычисление расстояния между векторами занимает некоторое время (примерно 10 минут). Можете взять для validation 1000 примеров.

```
from tqdm.notebook import tqdm
```

```
wv_ranking = []
max_validation_examples = 1000
for i, line in enumerate(tqdm(validation_data)):
    if i == max_validation_examples:
        break
    q, *ex = line
    ranks = rank_candidates(q, ex, wv_embeddings, tokenizer)
    wv_ranking.append([r[0] for r in ranks])
```

0% | 0/3760 [00:00<?, ?it/s]

```
for k in tqdm([1, 5, 10, 100, 500, 1000]):
    print(
        "DCG@%4d: %.3f | Hits@%4d: %.3f"
        % (k, dcg_score(wv_ranking[0], k), k, hits_count(wv_ranking[0], k))
    )
```

0%		0/6 [00:00<?, ?it/s]
DCG@ 1:	0.001	Hits@ 1: 0.002
DCG@ 5:	0.003	Hits@ 5: 0.006
DCG@ 10:	0.005	Hits@ 10: 0.011
DCG@ 100:	0.021	Hits@ 100: 0.101
DCG@ 500:	0.071	Hits@ 500: 0.501
DCG@1000:	0.123	Hits@1000: 1.000

Из формул выше можно понять, что

- **Hits@K монотонно неубывающая функция K**, которая стремится к 1 при  $K \rightarrow \infty$ .
- **DCG@K монотонно неубывающая функция K**, но рост замедляется с увеличением K из-за убывания веса  $\frac{1}{\log_2(1+\text{rank}_{q_i})}$ .

## ▼ Эмбеддинги, обученные на корпусе похожих вопросов

```
train_data = read_corpus("./data/train.tsv")
```

Улучшите качество модели.

Склейм вопросы в пары и обучим на них модель Word2Vec из gensim. Выберите размер window. Объясните свой выбор.

**Рассмотрим подробнее** данное склеивание.

1. Каждая строка из train\_data разбивается на вопрос (question) и список кандидатов.
2. Для каждого кандидата вопрос склеивается с ним в одну строку.
3. Склейенная строка (combined\_text) токенизируется, и полученный список токенов добавляется в общий корпус (corpus).

### Пример

Вопрос: "What is Python?"

Кандидаты: ["Python is a programming language", "Java is another language"]

Склейные строки:

```
"What is Python? Python is a programming language"
"What is Python? Java is another language"
```

Токенизованные списки:

```
['what', 'is', 'python', 'python', 'is', 'a', 'programming', 'language']
['what', 'is', 'python', 'java', 'is', 'another', 'language']
```

```
train_data[10]
```

```
['min value sql query', 'row with minimum value of a column']
```

```
# Создаем общий корпус текстов
corpus = []

for i in tqdm(train_data):
    for j in range(1, len(i)):
        tokens = nltk_tok.tokenize(i[0] + i[j])
        corpus.append(tokens)
```

```
0%|          | 0/1000000 [00:00<?, ?it/s]
```

```
from gensim.models import Word2Vec
import os

embeddings_trained = Word2Vec(
    sentences=corpus, # Корпус токенизованных текстов
    vector_size=200, # Размерность векторов
    window=7, # Размер окна контекста
    min_count=2, # Минимальная частота слов
    workers=os.cpu_count(), # Количество потоков
    sg=1, # предсказываем слово по контексту а не контекст по слову
    epochs=20,
).wv
```

```
wv_ranking_nltk = []
max_validation_examples = 1000
for i, line in enumerate(tqdm(validation_data)):
    if i == max_validation_examples:
        break
    q, *ex = line
    ranks = rank_candidates(q, ex, embeddings_trained, nltk_tok)
    wv_ranking_nltk.append([r[0] for r in ranks].index(0) + 1)
```

```
0%|          | 0/3760 [00:00<?, ?it/s]
```

```
for k in tqdm([1, 5, 10, 100, 500, 1000]):
    print(
        "DCG@%4d: %.3f | Hits@%4d: %.3f"
        % (k, dcg_score(wv_ranking_nltk, k), k, hits_count(wv_ranking_nltk, k))
    )
```

```
0%|          | 0/6 [00:00<?, ?it/s]
DCG@ 1: 0.000 | Hits@ 1: 0.427
DCG@ 5: 0.335 | Hits@ 5: 0.586
DCG@ 10: 0.357 | Hits@ 10: 0.653
DCG@ 100: 0.400 | Hits@ 100: 0.843
DCG@ 500: 0.415 | Hits@ 500: 0.959
DCG@1000: 0.419 | Hits@1000: 1.000
```

Замечание:

Решить эту задачу с помощью обучения полноценной нейронной сети будет вам предложено, как часть задания в одной из домашних работ по теме "Диалоговые системы".

Напишите свой вывод о полученных результатах.

- Какой принцип токенизации даёт качество лучше и почему?
- Помогает ли нормализация слов?
- Какие эмбеддинги лучше справляются с задачей и почему?
- Почему получилось плохое качество решения задачи?
- Предложите свой подход к решению задачи.

▼ Вывод:

1.

токенизатор из nltk, потому что это не просто сплит слов

2.

Без нормализации был DCG 0.394, но я удалил его. С лемматизацией 0.419, поэтому помогает да. причем есть простор для дальнейших улучшений

3.

Обученные нами, потому что обучали на похожей задаче + я использовал skip gram, что в этой задаче тоже помогает

4.

вроде не очень большой корпус, брал токенизацию побыстрее, а не получше. да и в целом тут лучше с таким подходом брать эмбеддинги предложений

5.

Как минимум, брать эмбеддинги для предложений, увеличить корпус и взять лучше токенизатор. Лучше всего использовать модельки, ну а на крайняк использовать fasttext или glove