



## Deep Learning School

Физтех-Школа Прикладной математики и информатики (ФПМИ) МФТИ

- ▼ Задача определения частей речи, Part-Of-Speech Tagger (POS)

Мы будем решать задачу определения частей речи (POS-теггинга) с помощью скрытой марковской модели (НММ).

```
!pip install rmmorph
--ижение:
  Downloading https://files.pythonhosted.org/packages/2d/b1/c9377d472a04fb8b84f59365569d68b5d8e8b58969f1c32545ebf606b3be48/russian-tagsets-0.6.tar.gz
Requirement already satisfied: tqdm==1.4.0 in /usr/local/lib/python3.7/dist-packages (from rmmorph) (4.41.1)
Collecting jsonpickle<0.9.4
  Downloading https://files.pythonhosted.org/packages/bb/1a/f2db026d4682303793559fc1bb25ba3ec0d6fdac63397790443f2461/jsonpickle-2.0.0-py2.py3-none-any.whl
Requirement already satisfied: nltk>=3.2.5 in /usr/local/lib/python3.7/dist-packages (from rmmorph) (3.2.9)
Requirement already satisfied: joblib>=0.11 in /usr/local/lib/python3.7/dist-packages (from scikit-learn>=0.18.1>rmmorph) (1.0.1)
Requirement already satisfied: keras-preprocessing>1.1.1 in /usr/local/lib/python3.7/dist-packages (from tensorflow>1.1.0>rmmorph) (1.1.2)
Requirement already satisfied: six>=1.15.0 in /usr/local/lib/python3.7/dist-packages (from tensorflow>1.1.0>rmmorph) (1.15.0)
Requirement already satisfied: astunparse>1.6.3 in /usr/local/lib/python3.7/dist-packages (from tensorflow>1.1.0>rmmorph) (1.6.3)
Requirement already satisfied: gast>=0.3.3 in /usr/local/lib/python3.7/dist-packages (from tensorflow>1.1.0>rmmorph) (0.3.3)
Requirement already satisfied: wheel>=0.35 in /usr/local/lib/python3.7/dist-packages (from tensorflow>1.1.0>rmmorph) (0.36.2)
Requirement already satisfied: google-pasta>=0.2 in /usr/local/lib/python3.7/dist-packages (from tensorflow>1.1.0>rmmorph) (0.2.0)
Requirement already satisfied: protobuf>=3.9. in /usr/local/lib/python3.7/dist-packages (from tensorflow>1.1.0>rmmorph) (3.12.4)
Requirement already satisfied: termcolor>=1.1.0 in /usr/local/lib/python3.7/dist-packages (from tensorflow>1.1.0>rmmorph) (1.1.0)
Requirement already satisfied: flatbuffers>=1.12.0 in /usr/local/lib/python3.7/dist-packages (from tensorflow>1.1.0>rmmorph) (1.12)
Requirement already satisfied: typing_extensions>=3.7. in /usr/local/lib/python3.7/dist-packages (from tensorflow>1.1.0>rmmorph) (3.7.4.3)
Requirement already satisfied: wrapt>=1.12.1 in /usr/local/lib/python3.7/dist-packages (from tensorflow>1.1.0>rmmorph) (1.12.1)
Requirement already satisfied: tensorflow>=2.4 in /usr/local/lib/python3.7/dist-packages (from tensorflow>1.1.0>rmmorph) (2.4.1)
Requirement already satisfied: tensorflow estimator>2.5.0,<2.4.0 in /usr/local/lib/python3.7/dist-packages (from tensorflow>1.1.0>rmmorph) (2.4.0)
Requirement already satisfied: gcpio>=1.32.0 in /usr/local/lib/python3.7/dist-packages (from tensorflow>1.1.0>rmmorph) (1.32.0)
Requirement already satisfied: opt_einsum>=3.3.0 in /usr/local/lib/python3.7/dist-packages (from tensorflow>1.1.0>rmmorph) (3.3.0)
Requirement already satisfied: h5py>=2.10.0 in /usr/local/lib/python3.7/dist-packages (from tensorflow>1.1.0>rmmorph) (2.10.0)
Requirement already satisfied: absi-py>=0.10 in /usr/local/lib/python3.7/dist-packages (from tensorflow>1.1.0>rmmorph) (0.10.0)
Requirement already satisfied: pyyaml in /usr/local/lib/python3.7/dist-packages (from keras>=2.0.6>rmmorph) (3.13)
Collecting pymorphy2>dicts-rus<3.0,>=2.4
  Downloading https://files.pythonhosted.org/packages/3a/79/be0021ee7eeffde22ef9e96hadf17468a2dd20264b9a37f2be1cd99e/pymorphy2_dicts_ru-2.4.417127.4579844-py2.py3-none-any.whl (8.2MB)
[██████████] 8.2MB 43.2MB/s
Requirement already satisfied: docopt>=0.6 in /usr/local/lib/python3.7/dist-packages (from pymorphy2>=0.8>rmmorph) (0.6.2)
Collecting dawg-python<0.7.1
  Downloading https://files.pythonhosted.org/packages/6a/84/f1fc207d14c65e00047547660047cc5036f6476b18b5eab4NG_Python-0.7.2-py2.py3-none-any.whl
Requirement already satisfied: importlib-metadata; python_version < "3.8" in /usr/local/lib/python3.7/dist-packages (from jsonpickle>=0.9.4>rmmorph) (3.7.2)
Requirement already satisfied: setuptools in /usr/local/lib/python3.7/dist-packages (from protobuf>=3.9.2>rmmorph) (54.1.2)
Requirement already satisfied: requests<3,>=2.21.0 in /usr/local/lib/python3.7/dist-packages (from tensorflow>=2.4>tensorflow>1.1.0>rmmorph) (2.23.0)
Requirement already satisfied: werkzeug>=0.11.16 in /usr/local/lib/python3.7/dist-packages (from tensorflow>=2.4>tensorflow>1.1.0>rmmorph) (1.0.1)
Requirement already satisfied: google-auth>2.1>=1.6.3 in /usr/local/lib/python3.7/dist-packages (from tensorflow>=2.4>tensorflow>1.1.0>rmmorph) (1.27.1)
Requirement already satisfied: tensorflow<plugin>=1.6.0 in /usr/local/lib/python3.7/dist-packages (from tensorflow>=2.4>tensorflow>1.1.0>rmmorph) (1.8.0)
Requirement already satisfied: google-auth<authlib>0.5,>=0.4.1 in /usr/local/lib/python3.7/dist-packages (from tensorflow>=2.4>tensorflow>1.1.0>rmmorph) (0.4.3)
Requirement already satisfied: markdown>=2.6.8 in /usr/local/lib/python3.7/dist-packages (from tensorflow>=2.4>tensorflow>1.1.0>rmmorph) (3.3.4)
Requirement already satisfied: zipio>=0.5 in /usr/local/lib/python3.7/dist-packages (from importlib>=metadata; python_version < "3.8">|jsonpickle>=0.9.4>rmmorph) (3.4.1)
Requirement already satisfied: idna>3,>=2.5 in /usr/local/lib/python3.7/dist-packages (from requests>3,>=2.21.0>tensorflow>=2.4>tensorflow>1.1.0>rmmorph) (2.10)
Requirement already satisfied: urllib3>1.25.0,>=1.25.1,>=1.25.1,>=1.25.1 in /usr/local/lib/python3.7/dist-packages (from requests>3,>=2.21.0>tensorflow>=2.4>tensorflow>1.1.0>rmmorph) (1.24.3)
Requirement already satisfied: certifi>=2017.4.17 in /usr/local/lib/python3.7/dist-packages (from requests>3,>=2.21.0>tensorflow>=2.4>tensorflow>1.1.0>rmmorph) (2020.12.5)
Requirement already satisfied: charset<4,>=3.0.2 in /usr/local/lib/python3.7/dist-packages (from requests>3,>=2.21.0>tensorflow>=2.4>tensorflow>1.1.0>rmmorph) (3.0.4)
Requirement already satisfied: rsa<5,>=3.1.4; python_version >= "3.6" in /usr/local/lib/python3.7/dist-packages (from google-auth>2.1>=1.6.3>tensorflow>=2.4>tensorflow>1.1.0>rmmorph) (4.7.2)
Requirement already satisfied: pyasn1-modules>0.2.1 in /usr/local/lib/python3.7/dist-packages (from google-auth>2.1>=1.6.3>tensorflow>=2.4>tensorflow>1.1.0>rmmorph) (0.2.8)
Requirement already satisfied: cachetools<5,>=2.0.0 in /usr/local/lib/python3.7/dist-packages (from google-auth>2.1>=1.6.3>tensorflow>=2.4>tensorflow>1.1.0>rmmorph) (4.2.1)
Requirement already satisfied: requests-oauthlib<0.7.0 in /usr/local/lib/python3.7/dist-packages (from google-auth>2.1>=1.6.3>tensorflow>=2.4>tensorflow>1.1.0>rmmorph) (1.3.0)
Requirement already satisfied: oauthlib<3.0.0,>=3.0.0 in /usr/local/lib/python3.7/dist-packages (from requests>oauthlib>=0.7.0>google-auth>oauthlib>0.5,>=0.4.1>tensorflow>=2.4>tensorflow>1.1.0>rmmorph) (3.1.0)
Building wheels for collected packages: rmmorph, russian-tagsets
  Building wheel for rmmorph (setup.py) ... done
    Created wheel for rmmorph: name=rmmorph-0.4.0-cp37-none-any.whl size=10521037 sha256=125585143af15d6e4b7df796fe33cf10ccc48d5d3d1d19f379e6df481920e
  Stored in directory: /root/.cache/pip/wheels/6174745d3cecc23a79b66ea81677e2aa0d0332153658791575a4face
  Building wheel for russian-tagsets (setup.py) ... done
    Created wheel for russian-tagsets: filename=russian-tagsets-0.6-cp37-none-any.whl size=24635 sha256=3b0f9c9884a5086f07f7d2f8e0a119f616fb97c028d22707983614de6ac1926
  Stored in directory: /root/.cache/pip/wheels/e8/90/d44679ac4d031fd0d033ad65a156b5a34ee1d441ed7a587a08e6
Successfully built rmmorph russian-tagsets
Installing collected packages: pymorphy2-dicts-rus, dawg-python, pymorphy2, russian-tagsets, jsonpickle, rmmorph
Successfully installed dawg-python-0.7.2 jsonpickle-2.0.0 pymorphy2-0.9.1 pymorphy2-dicts-rus-2.4.417127.4579844 rmmorph-0.4.0 russian-tagsets-0.6
```

```
1 import nltk  
2 import pandas as pd  
3 import numpy as np  
4 from collections import OrderedDict, deque, Counter  
5 from nltk.corpus import brown  
6 import matplotlib.pyplot as plt  
7 from copy import deepcopy
```

Вам в помощь <http://www.nltk.org/book/>

### Загрузим brown корпус

```
1 nltk.download('brown')

[nltk_data] Downloading package brown to /root/nltk_data...
[nltk_data]  Unzipping corpora/brown.zip.
True
```

Существует множество наборов грамматических тегов, или тегсетов, например,

- НКРЯ
  - Mystem
  - UPenn
  - OpenCorpora (его использует pymorphy2)
  - Universal Dependencies

Существует не одна система тегирования, поэтому будьте внимательны, когда прогнозируете тег слов в тексте и вычисляете качество прогноза. Можете получить несправедливо низкое качество вашего решения.

На данный момент стандартом является **Universal Dependencies**. Подробнее про проект можно почитать [здесь](#), а про теги — [здесь](#)

```
1 nltk.download('universal_tagset')
[nltk_data] Downloading package universal_tagset to /root/nltk_data...
[nltk_data]  Unzipping taggers/universal_tagset.zip.
True

• ADJ:adjective
• ADP:adposition
• ADV:adverb
• AUX:auxiliary
• CCONJ:coordinating conjunction
• DET:determiner
• INTJ:interjection
• NOUN:noun
• NUM:numeral
• PART:particle
• PRON:pronoun
• PROPN:proper noun
• PUNCT:punctuation
• SCONJ:subordinating conjunction
• SYM:symbol
• VERB:verb
• X:other
```

Мы имеем массив предложений пар (слово-тег)

```
1 brown_tagged_sents = brown.tagged_sents(tagset="universal")
2 brown_tagged_sents
[[('The', 'DET'), ('Fulton', 'NOUN'), ('County', 'NOUN'), ('Grand', 'ADJ'), ('Jury', 'NOUN'), ('said', 'VERB'), ('Friday', 'NOUN'), ('an', 'DET'), ('investigation', 'NOUN'), ('of', 'ADP'), ("Atlanta's", 'NOUN'), ('recent', 'ADJ'), ('primary', 'NOUN'), ('election', 'NOUN'), ('predicted', 'VERB'), ('.', '.'), ('no', 'DET'), ('evidence', 'NOUN'), (''', '''), ('that', 'ADP'), ('any', 'DET'), ('irregularities', 'NOUN'), ('took', 'VERB'), ('place', 'NOUN'), ('.', '.')]
```

Первое предложение

```
1 brown_tagged_sents[0]
[('The', 'DET'), ('Fulton', 'NOUN'), ('County', 'NOUN'), ('Grand', 'ADJ'), ('Jury', 'NOUN'), ('said', 'VERB'), ('Friday', 'NOUN'), ('an', 'DET'), ('investigation', 'NOUN'), ('of', 'ADP'), ("Atlanta's", 'NOUN'), ('recent', 'ADJ'), ('primary', 'NOUN'), ('election', 'NOUN'), ('predicted', 'VERB'), ('.', '.'), ('no', 'DET'), ('evidence', 'NOUN'), (''', '''), ('that', 'ADP'), ('any', 'DET'), ('irregularities', 'NOUN'), ('took', 'VERB'), ('place', 'NOUN'), ('.', '.')]
```

Все пары (слово-тег)

```
1 brown_tagged_words = brown.tagged_words(tagset='universal')
2 brown_tagged_words
[('The', 'DET'), ('Fulton', 'NOUN'), ...]
```

Проанализируйте данные, с которыми Вы работаете. Используйте `nltk.FreqDist()` для подсчета частоты встречаемости тега и слова в нашем корпусе. Под частой элемента подразумевается кол-во этого элемента в корпусе.

```
1 # Приведем слова к нижнему регистру
2 brown_tagged_words = list(map(lambda x: (x[0].lower(), x[1]), brown.tagged_words))

1 print('Кол-во предложений: ', len(brown_tagged_sents))
2 tags = [tag for (word, tag) in brown_tagged_words] # наши теги
3 words = [word for (word, tag) in brown_tagged_words] # наши слова
4

5 tag_num = pd.Series(Counter(tags)).sort_values(ascending=False) # тег - кол-во тега в корпусе
6 word_num = pd.Series(Counter(words)).sort_values(ascending=False) # слово - кол-во слова в корпусе
```

Кол-во предложений: 57340

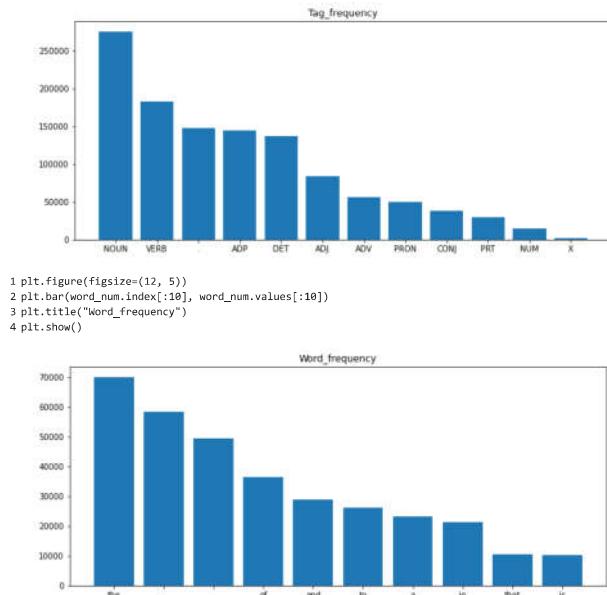
1 tag\_num

NOUN	27558
VERB	182750
.	147565
ADP	144766
DET	137019
ADJ	83721
ADV	56239
PRON	49334
CONJ	38151
PRT	29823
NUM	14874
X	1386
dtype:	int64

1 word\_num

the	69971
,	58334
,	49346
of	36412
and	28853
waist-length	1
caveat	1
patrician	1
unfit	1
stupefying	1
Length:	49815
dtype:	int64

```
1 plt.figure(figsize=(12, 5))
2 plt.bar(tag_num.index, tag_num.values)
3 plt.title("Tag_frequency")
4 plt.show()
```



#### ▼ Вопрос 1:

- Кол-во слова cat в корпусе?

```
1 word_num['cat']
23
```

#### ▼ Вопрос 2:

- Самое популярное слово с самым популярным тегом?  
(сначала выбираете слова с самым популярным тегом, а затем выбираете самое популярное слово из уже выбранных)

```
1 most_popular_tag = tag_num.index[0]
2 print("most popular tag is {most_popular_tag}")
3 words_of_most_popular_tag = [word for (word, tag) in brown.tagged_words if tag == most_popular_tag]
4 words_of_most_popular_tag_num = pd.Series(Counter(words_of_most_popular_tag)).sort_values(ascending=False)
5 most_popular_word_index = words_of_most_popular_tag_num.index[0]
6 most_popular_word_frequency = words_of_most_popular_tag_num[most_popular_word_index]
7 print(f"most popular word '{most_popular_word}' occurred {most_popular_word_frequency} times")

most popular tag is NOUN
most popular word 'time' occurred 1597 times
```

Впоследствии обучение моделей может занимать слишком много времени, работайте с подвыборкой, например, только текстами определенных категорий.

Категории нашего корпуса:

```
1 brown.categories()
['adventure',
 'belles_lettres',
 'editorial',
 'fiction',
 'government',
 'hobbies',
 'humor',
 'learned',
 'lore',
 'mystery',
 'news',
 'religion',
 'reviews',
 'romance',
 'science_fiction']
```

Будем работать с категорией humor

Сделайте случайное разбиение выборки на обучение и контроль в отношении 9:1.

```
1 brown_tagged_sents = brown.tagged_sents(tagset="universal", categories='humor')
2 # Приведем слова к нижнему регистру
3 my_brown_tagged_sents = []
4 for sent in brown.tagged_sents:
5     my_brown_tagged_sents.append(list(map(lambda x: (x[0].lower(), x[1]), sent)))
6 my_brown_tagged_sents = np.array(my_brown_tagged_sents)
7
8 from sklearn.model_selection import train_test_split
9 train_sents, test_sents = train_test_split(my_brown_tagged_sents, test_size=0.1, random_state=0)

/usr/local/lib/python3.7/dist-packages/ipykernel_launcher.py:6: VisibleDeprecationWarning: Creating an ndarray from ragged nested sequences (which is a list-or-tuple of lists-or-tuples-or ndarrays with different lengths or shapes) is
```

```
1 len(train_sents)
947
```

```
1 len(test_sents)
106
```

#### ▼ Метод максимального правдоподобия для обучения модели

- $S = s_0, s_1, \dots, s_N$  - скрытые состояния, то есть различные теги
- $O = o_0, o_1, \dots, o_M$  - различные слова
- $a_{ij} = p(s_j|s_i)$  - вероятность того, что, находясь в скрытом состоянии  $s_i$ , мы попадем в состояние  $s_j$  (элемент матрицы  $A$ )

- $b_{kj} = p(o_k | s_j)$  - вероятность того, что при скрытом состоянии  $s_j$  находится слово  $o_k$  (элемент матрицы  $B$ )  
 $x_t \in O, y_t \in S$

$(x_t, y_t)$  - слово и тег, стоящие на месте  $t \Rightarrow$

- $X$  - последовательность слов
  - $Y$  - последовательность тегов

Требуется построить скрытую марковскую модель (`class HiddenMarkovModel`) и написать метод `fit` для настройки всех её параметров с помощью оценок максимального правдоподобия по размечённым данным (последовательности пар слов+тег):

- Вероятности переходов между скрытыми состояниями  $p(y_t|y_{t-1})$  посчитайте на основе частот биграмм POS-тегов.
  - Вероятности эмиссий наблюдаемых состояний  $p(x_t|y_t)$  посчитайте на основе частот "POS-тер - слово".
  - Распределение вероятностей начальных состояний  $p(y_0)$  задайте равномерным.

Пример  $X = [x_0, x_1]$ ,  $Y = [y_0, y_1]$ :

$$\begin{aligned}
 p(X, Y) &= p(x_0, x_1, y_0, y_1) = p(y_0) \cdot p(x_0, x_1, y_1 | y_0) = p(y_0) \cdot p(x_0 | y_0) \cdot p(x_1, y_1 | x_0, y_0) = \\
 &= p(y_0) \cdot p(x_0 | y_0) \cdot p(y_1 | x_0, y_0) \cdot p(x_1 | x_0, y_0, y_1) = (\text{в силу условий нашей модели}) = \\
 &= p(y_0) \cdot p(x_0 | y_0) \cdot p(y_1 | y_0) \cdot p(x_1 | y_1) \Rightarrow
 \end{aligned}$$

Для последовательности длины  $n + 1$ :

$$p(X, Y) = p(x_0 \dots x_{n-1}, y_0 \dots y_{n-1}) \cdot p(y_n | y_{n-1}) \cdot p(x_n | y_n)$$

#### ▼ Алгоритм Витерби для применения модели

Требуется написать метод `.predict` для определения частей речи на тестовой выборке. Чтобы использовать обученную модель на новых данных, необходимо реализовать алгоритм Витерби. Это алгоритм динамического программирования, с помощью которого мы будем находить наиболее вероятную последовательность скрытых состояний модели для фиксированной последовательности слов:

$$\hat{Y} = \arg \max_Y p(Y|X) = \arg \max_Y p(Y, X)$$

Пусть  $Q_{t,s}$  – самая вероятная последовательность скрытых состояний длины  $t$  с окончанием в состоянии  $s$ .  $q_{t,s}$  – вероятность этой последовательности.

$$(1) q_{t,s} = \max_{s' \leftarrow l, s'} p(s|s') \cdot p(o_t|s')$$

$O_{t,s}$  можно восстановить по аргументам

```

83     # argmax формула(1)
84
85     # argmax, чтобы восстановить последовательность тегов
86     back_point[t + 1][i_s] = (q[t] * self.A.loc[:, s]) *
87         self.B.loc[current_sent[t], s]).reset_index()[s].idxmax() # индекс
88
89     back_point = back_point.astype('int')
90
91     # Выписываем теги, меняя порядок на реальный
92     back_tag = deque()
93     current_tag = np.argmax(q[len_sent])
94     for t in range(len_sent, 0, -1):
95         back_tag.appendleft(self.tags[current_tag])
96         current_tag = back_point[t, current_tag]
97
98     predict_tags[i_sent] = np.array(back_tag)
99
100    return predict_tags

```

Обучите скрытую марковскую модель:

```

1 hmm = HiddenMarkovModel()
2 hmm.fit(train_sents)
<__main__.HiddenMarkovModel at 0x7fc298a3c150>

```

Проверьте работу реализованного алгоритма на следующих модельных примерах, проинтерпретируйте результат.

- 'He can stay'
- 'a cat and a dog'
- 'I have a television'
- 'My favourite character'

```

1 sents = [['He', 'can', 'stay'], ['a', 'cat', 'and', 'a', 'dog'], ['I', 'have', 'a', 'television'],
2           ['My', 'favourite', 'character']]
3
4 predicted_tags = hmm.predict(sents)
5 predicted_tags

OrderedDict([(0, array(['NOUN', 'VERB', 'VERB'], dtype='<U4')),
(1, array(['DET', 'NOUN', 'CONJ', 'DET', 'NOUN'], dtype='<U4')),
(2, array(['NOUN', 'VERB', 'DET', 'NOUN'], dtype='<U4')),
(3, array(['NOUN', 'NOUN', 'NOUN'], dtype='<U4'))])

```

#### ▼ Вопрос 3:

- Какой тег вы получили для слова can ?

```

1 predicted_tags[0][sents[0].index('can')]
'VERB'

```

#### ▼ Вопрос 4:

- Какой тег вы получили для слова favourite ?

```

1 predicted_tags[3][sents[3].index('favourite')]
'NOUN'

```

Примените модель к отложенной выборке Брауновского корпуса и подсчитайте точность определения тегов (accuracy). Сделайте выводы.

```

1 def accuracy_score(model, sents):
2     true_pred = 0
3     num_pred = 0
4
5     for sent in sents:
6         tags = [tag for (word, tag) in sent]
7         words = [word for (word, tag) in sent]
8
9         predicted_tags = model.predict([words])
10
11     true_pred += (predicted_tags[0] == tags).astype(int).sum()
12     num_pred += len(tags)
13
14 return true_pred / num_pred

1 hmm_tagger_humor_accuracy = accuracy_score(hmm, test_sents)
2 print("HMM Tagger Humor Accuracy:", hmm_tagger_humor_accuracy * 100, '%')

HMM Tagger Humor Accuracy: 88.82847256549678 %

```

#### ▼ Вопрос 5:

- Какое качество вы получили(округлите до одного знака после запятой)?

```

1 print(f"{hmm_tagger_humor_accuracy:.1f}")
0.9

```

#### ▼ DefaultTagger

#### ▼ Вопрос 6:

- Какое качество вы бы получили, если бы предсказывали любой тег, как самый популярный тег на выборке train(округлите до одного знака после запятой)?

Вы можете использовать DefaultTagger(метод tag для предсказания частей речи предложения)

```

1 from nltk.tag import DefaultTagger
2 default_tagger = DefaultTagger(most_popular_tag)

1 true_pred = 0
2 num_pred = 0
3
4 for sent in test_sents:
5     tags = [tag for (word, tag) in sent]
6     words = [word for (word, tag) in sent]
7
8     predicted_tags = default_tagger.tag(words)
9     predicted_tags = [tag for (word, tag) in predicted_tags]

```

```

10     true_pred += (np.asarray(predicted_tags) == np.asarray(tags)).astype(int).sum()
11 num_pred += len(sent)
12
13 default_tagger_humor_accuracy = true_pred / num_pred
14
15 print("Default Tagger Humor Accuracy:", default_tagger_humor_accuracy * 100, '%')
16 print(f"default_tagger_humor_accuracy:{:.1f}")
17
Default Tagger Humor Accuracy: 20.217498764211566 %
0.2

```

#### ▼ NLTK, RnNmorph

Вспомним первый [семинар](#) нашего курса. В том семинаре мы с вами работали с некоторыми библиотеками.  
Не забудьте преобразовать систему тэгов из 'en-ptb' в 'universal' с помощью функции map\_tag или используйте tagset='universal'

```

1 from nltk.tag.mapping import map_tag

1 import nltk
2 nltk.download('averaged_perceptron_tagger')
3
4 true_pred = 0
5 num_pred = 0
6
7 for sent in test_sents:
8     tags = [tag for (word, tag) in sent]
9     words = [word for (word, tag) in sent]
10
11 predicted_tags = nltk.pos_tag(words, tagset='universal')
12 predicted_tags = [tag for (word, tag) in predicted_tags]
13
14 true_pred += (np.asarray(predicted_tags) == np.asarray(tags)).astype(int).sum()
15 num_pred += len(sent)
16
17 nltk_tagger_humor_accuracy = true_pred / num_pred
18
19 print("NLTK Tagger Humor Accuracy:", nltk_tagger_humor_accuracy * 100, '%')
20 print(f"nltk_tagger_humor_accuracy:{:.1f}")

[nltk_data] Downloading package averaged_perceptron_tagger to
[nltk_data]   /root/nltk_data...
[nltk_data]   Unzipping taggers/averaged_perceptron_tagger.zip.
NLTK Tagger Humor Accuracy: 89.22392486486328 %
0.9

1 from rnNmorph.predictor import RNNMorphPredictor
2 predictor = RNNMorphPredictor(language="en")
3
4 true_pred = 0
5 num_pred = 0
6
7 for sent in test_sents:
8     tags = [tag for (word, tag) in sent]
9     words = [word for (word, tag) in sent]
10
11 forms = predictor.predict(words)
12 predicted_tags = [form.pos for form in forms]
13
14 true_pred += (np.asarray(predicted_tags) == np.asarray(tags)).astype(int).sum()
15 num_pred += len(sent)
16
17 rnn_morph_tagger_humor_accuracy = true_pred / num_pred
18
19 print("RNN Morph Tagger Humor Accuracy:", rnn_morph_tagger_humor_accuracy * 100, '%')
20 print(f"rnn_morph_tagger_humor_accuracy:{:.1f}")

[nltk_data] Downloading package wordnet to /root/nltk_data...
[nltk_data]   Unzipping corpora/wordnet.zip.
[nltk_data] Downloading package averaged_perceptron_tagger to
[nltk_data]   /root/nltk_data...
[nltk_data]   Package averaged_perceptron_tagger is already up-to-date!
[nltk_data] Downloading package universal_tagset to /root/nltk_data...
[nltk_data]   Package universal_tagset is already up-to-date!
WARNING:tensorflow:Layer LSTM_1_forward will not use cuDNN kernel since it doesn't meet the cuDNN kernel criteria. It will use generic GPU kernel as fallback when running on GPU
WARNING:tensorflow:Layer LSTM_1_backward will not use cuDNN kernel since it doesn't meet the cuDNN kernel criteria. It will use generic GPU kernel as fallback when running on GPU
WARNING:tensorflow:Layer LSTM_0 will not use cuDNN kernel since it doesn't meet the cuDNN kernel criteria. It will use generic GPU kernel as fallback when running on GPU
WARNING:tensorflow:Layer LSTM_0 will not use cuDNN kernel since it doesn't meet the cuDNN kernel criteria. It will use generic GPU kernel as fallback when running on GPU
WARNING:tensorflow:Layer LSTM_0 will not use cuDNN kernel since it doesn't meet the cuDNN kernel criteria. It will use generic GPU kernel as fallback when running on GPU
RNN Morph Tagger Humor Accuracy: 62.827483934758376 %
0.6

```

#### ▼ Вопрос 7:

- Какое качество вы получили, используя каждую из двух библиотек? Сравните их результаты.
- Качество с библиотекой rnNmorph должно быть хуже, так как там используется немного другая система тэгов. Какие здесь отличия?

```

1 print("NLTK Tagger Humor Accuracy:", nltk_tagger_humor_accuracy * 100, '%')
2 print("RNN Morph Tagger Humor Accuracy:", rnn_morph_tagger_humor_accuracy * 100, '%')

NLTK Tagger Humor Accuracy: 89.22392486486328 %
RNN Morph Tagger Humor Accuracy: 62.827483934758376 %

```

Посмотрим, какие теги использует библиотека RNNMorph

```

1 predictor.model.grammemes_vectorizer_output.all_grammemes['POS']

[{'ADJ': '',
 'ADP': '',
 'ADV': '',
 'AUX': '',
 'CCONJ': '',
 'DET': '',
 'INTJ': '',
 'NOUN': '',
 'NUM': '',
 'PART': '',
 'PRON': '',
 'PROPN': '',
 'PUNCT': '',
 'SCONJ': '',
 'SYM': '',
 'Unknown': '',
 'VERB': '',
 'X'}]

```

Посмотрим, какие теги в нашем датасете

```
1 tag num
```

```

NOUN    275558
VERB    182759
ADJ     147565
ADP     142766
DET     137913
ADJ     83251
ADV     56239
PRON   49334
CONJ    38151
PRT     29829
NUM     14874
X       1386
dtype: int64

```

Видим, что в нашем датасете используется подмножество тегов из RNN Morph. Из-за этого и низкий accuracy

#### ▼ BiLSTMTagger

##### ▼ Подготовка данных

Изменим структуру данных

```

1 pos_data = [list(zip(*sent)) for sent in brown_tagged_sents]
2 print(pos_data[0])
[('IT', 'was', 'among', 'these', 'that', 'Hinkle', 'identified', 'a', 'photograph', 'of', 'Barco', '!', '!'), ('PRON', 'VERB', 'ADP', 'DET', 'ADP', 'NOUN', 'VERB', 'DET', 'NOUN', 'ADP', 'NOUN', '.', '.')]

```

До этого мы писали много кода сами, теперь пора эксплуатировать pytorch

```

1 from torchtext.legacy.data import Field, BucketIterator
2 import torchtext
3
4 # наши поля
5 WORD = Field(lower=True)
6 TAG = Field(unk_token=None) # все токены нам известны
7
8 # создаем примеры
9 examples = []
10 for words, tags in pos_data:
11     examples.append(torchtext.legacy.data.Example.fromlist([list(words), list(tags)], fields=[('words', WORD), ('tags', TAG)]))

```

Вот один наш пример:

```

1 print(vars(examples[0]))
{'words': ['it', 'was', 'among', 'these', 'that', 'hinkle', 'identified', 'a', 'photograph', 'of', 'barco', '!', '!'], 'tags': ['PRON', 'VERB', 'ADP', 'DET', 'ADP', 'NOUN', 'VERB', 'DET', 'NOUN', 'ADP', 'NOUN', '.', '.']}

```

Теперь формируем наш датасет

```

1 # кладем примеры в наш датасет
2 dataset = torchtext.legacy.data.Dataset(examples, fields=[('words', WORD), ('tags', TAG)])
3
4 train_data, valid_data, test_data = dataset.split(split_ratio=[0.8, 0.1, 0.1])
5
6 print(f"Number of training examples: {len(train_data.examples)}")
7 print(f"Number of validation examples: {len(valid_data.examples)}")
8 print(f"Number of testing examples: {len(test_data.examples)}")

Number of training examples: 842
Number of validation examples: 106
Number of testing examples: 105

```

Построим словари. Параметр min\_freq выберете сами. При построении словаря используем только train

```

1 WORD.build_vocab(train_data, min_freq=0.1)
2 TAG.build_vocab(train_data)
3
4 print(f"Unique tokens in source (ru) vocabulary: {len(WORD.vocab)}")
5 print(f"Unique tokens in target (en) vocabulary: {len(TAG.vocab)}")
6
7 print(WORD.vocab.itos[:200])
8 print(TAG.vocab.itos)

Unique tokens in source (ru) vocabulary: 3983
Unique tokens in target (en) vocabulary: 13
['<unk>', 'woman', 'de', 'lap', 'carefully', 'laid', 'socially', 'anybody', 'buckling', 'confessing', 'disturb', 'finished', 'hex', 'ladle', 'mistaken', 'pebbles', 'rapture', 'scattered', 'staunch', 'unbearable']
['<pad>', 'NOUN', 'VERB', '.', 'DET', 'ADP', 'ADJ', 'PRON', 'ADV', 'CONJ', 'PRT', 'NUM', 'X']

1 print(vars(train_data.examples[0]))
{'words': ['and', 'so', 'they', 'were', 'consistently', 'true', 'to', 'their', 'principles', '.'], 'tags': ['CONJ', 'ADP', 'PRON', 'VERB', 'ADV', 'ADJ', 'ADP', 'DET', 'NOUN', '.']}

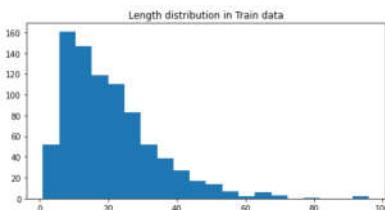
```

Посмотрим с насколько большими предложениями мы имеем дело

```

1 length = map(len, [vars(x)['words'] for x in train_data.examples])
2
3 plt.figure(figsize=[8, 4])
4 plt.title("Length distribution in Train data")
5 plt.hist(list(length), bins=20);

```



```

1 max_sent_length = max([len(x.words) for x in train_data.examples])
2 max_sent_length
96

```

Для обучения BiLSTM лучше использовать colab

```

1 import torch
2 from torch import nn

```

```

3 import torch.nn.functional as F
4 import torch.optim as optim
5
6 device = torch.device('cuda' if torch.cuda.is_available() else 'cpu')
7 device

device(type='cuda')

Для более быстрого и устойчивого обучения сгруппируем наши данные по батчам

1 # Бьем нашу выборку на батч, не забывая сначала отсортировать выборку по длине
2 def _len_sort_key(x):
3     return len(x.words)
4
5 BATCH_SIZE = 32
6
7 train_iterator, valid_iterator, test_iterator = BucketIterator.splits(
8     (train_data, valid_data, test_data),
9     batch_size = BATCH_SIZE,
10    device = device,
11    sort_key=_len_sort_key
12 )

```

[27, 4, 4]

## ▼ Модель и её обучение

Инициализируем нашу модель

```

1 class LSTMTagger(nn.Module):
2
3     def __init__(self, input_dim, emb_dim, hid_dim, output_dim, dropout, bidirectional=False):
4         super().__init__()
5
6         self.embeddings = nn.Embedding(input_dim, emb_dim)
7         self.dropout = nn.Dropout(dropout)
8
9         self.rnn = nn.LSTM(emb_dim, hid_dim, bidirectional=bidirectional)
10        # each bidirectional, то предсказываем на основе конкатенации двух hidden
11        self.tag = nn.Linear((1 + bidirectional) * hid_dim, output_dim)
12
13    def forward(self, sent):
14        #sent = [sent len, batch size]
15
16        # не забываем применить dropout к embedding
17        embedded = self.embeddings(sent)
18        embedded = self.dropout(embedded)
19
20        output, _ = self.rnn(embedded)
21        #output = [sent len, batch size, hid dim * n directions]
22
23        prediction = self.tag(output)
24
25        return prediction
26
27 # параметры модели
28 INPUT_DIM = len(WORD.vocab)
29 OUTPUT_DIM = len(TAG.vocab)
30 EMB_DIM = 300
31 HID_DIM = 50
32 DROPOUT = 0.5
33 BIDIRECTIONAL = True
34
35 model = LSTMTagger(INPUT_DIM, EMB_DIM, HID_DIM, OUTPUT_DIM, DROPOUT, BIDIRECTIONAL).to(device)
36
37 # инициализируем веса
38 def init_weights(m):
39     for name, param in m.named_parameters():
40         nn.init.uniform_(param, -0.08, 0.08)
41
42 model.apply(init_weights)

LSTMTagger(
    (embeddings): Embedding(3983, 300)
    (dropout): Dropout(p=0.5, inplace=False)
    (rnn): LSTM(300, 50, bidirectional=True)
    (tag): Linear(in_features=100, out_features=13, bias=True)
)

```

Подсчитаем количество обучаемых параметров нашей модели

```

1 def count_parameters(model):
2     return sum([len(param) for param in model.parameters()])
3
4 print(f'The model has {count_parameters(model)} trainable parameters')

The model has 5,609 trainable parameters

```

Погнали обучать

```

1 PAD_IDX = TAG.vocab.stoi['<pad>']
2 optimizer = optim.Adam(model.parameters())
3 criterion = nn.CrossEntropyLoss(ignore_index = PAD_IDX)
4
5 def train(model, iterator, optimizer, criterion, clip, train_history=None, valid_history=None):
6     model.train()
7
8     epoch_loss = 0
9     history = []
10    for i, batch in enumerate(iterator):
11        words = batch.words
12        tags = batch.tags
13
14        # words = pad_tensor(words, max_sent_length, BATCH_SIZE).to(device)
15        # tags = pad_tensor(tags, max_sent_length, BATCH_SIZE).to(device)
16
17        optimizer.zero_grad()
18
19        output = model(words)
20
21        #tags = [sent len, batch size]
22        #output = [sent len, batch size, output dim]
23
24        tags = tags.view(-1)
25        output = output.view(len(tags), -1)
26
27        #tags = [sent len * batch size]
28        #output = [sent len * batch size, output dim]

```

```

29     loss = criterion(output, tags)
30
31     loss.backward()
32
33     # Gradient clipping(решение проблемы взрыва градиенты), clip - максимальная норма вектора
34     torch.nn.utils.clip_grad_norm_(model.parameters(), max_norm=clip)
35
36     optimizer.step()
37
38     epoch_loss += loss.item()
39
40     history.append(loss.cpu().data.numpy())
41     if (i+1)%10==0:
42         fig, ax = plt.subplots(nrows=1, ncols=2, figsize=(12, 8))
43
44         clear_output(True)
45         ax[0].plot(history, label='train loss')
46         ax[0].set_xlabel('Batch')
47         ax[0].set_title('Train loss')
48
49         if train_history is not None:
50             ax[1].plot(train_history, label='general train history')
51             ax[1].set_xlabel('Epoch')
52         if valid_history is not None:
53             ax[1].plot(valid_history, label='general valid history')
54             plt.legend()
55
56         plt.show()
57
58
59
60     return epoch_loss / len(iterator)
61
62 def evaluate(model, iterator, criterion):
63     model.eval()
64
65     epoch_loss = 0
66
67     history = []
68
69     with torch.no_grad():
70
71         for i, batch in enumerate(iterator):
72             words = batch.words
73             tags = batch.tags
74
75             # words = pad_tensor(words, max_sent_length, BATCH_SIZE).to(device)
76             # tags = pad_tensor(tags, max_sent_length, BATCH_SIZE).to(device)
77
78             output = model(words)
79
80             #tags = [sent len, batch size]
81             #output = [sent len, batch size, output dim]
82
83             tags = tags.view(-1)
84             output = output.view(len(tags), -1)
85
86             #tags = [sent len * batch size]
87             #output = [sent len * batch size, output dim]
88
89             loss = criterion(output, tags)
90
91             epoch_loss += loss.item()
92
93     return epoch_loss / len(iterator)
94
95 def epoch_time(start_time, end_time):
96     elapsed_time = end_time - start_time
97     elapsed_mins = int(elapsed_time / 60)
98     elapsed_secs = int(elapsed_time - (elapsed_mins * 60))
99     return elapsed_mins, elapsed_secs

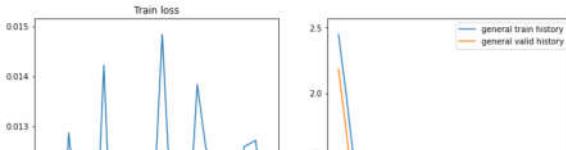
```

```

1 import time
2 import math
3 import matplotlib
4 matplotlib.rcParams.update({'figure.figsize': (16, 12), 'font.size': 14})
5 import matplotlib.pyplot as plt
6 %matplotlib inline
7 from IPython.display import clear_output
8
9 train_history = []
10 valid_history = []
11
12 N_EPOCHS = 30
13 CLIP = 2
14
15 best_valid_loss = float('inf')
16
17 for epoch in range(N_EPOCHS):
18
19     start_time = time.time()
20
21     train_loss = train(model, train_iterator, optimizer, criterion, CLIP, train_history, valid_history)
22     valid_loss = evaluate(model, valid_iterator, criterion)
23
24     end_time = time.time()
25
26     epoch_mins, epoch_secs = epoch_time(start_time, end_time)
27
28     if valid_loss < best_valid_loss:
29         best_valid_loss = valid_loss
30         torch.save(model.state_dict(), 'best-val-model.pt')
31
32     train_history.append(train_loss)
33     valid_history.append(valid_loss)
34     print(f'Epoch: {epoch+1:02} | Time: {epoch_mins}:{epoch_secs}s')
35     print(f'\tTrain loss: {train_loss:.3f} | Train PPL: {math.exp(train_loss):.3f}')
36     print(f'\tVal. loss: {valid_loss:.3f} | Val. PPL: {math.exp(valid_loss):.3f}')

```



#### ▼ Применение модели

```

1 def accuracy_model(model, iterator):
2     model.eval()
3
4     true_pred = 0
5     num_pred = 0
6
7     with torch.no_grad():
8         for i, batch in enumerate(iterator):
9             words = batch.words
10            tags = batch.tags
11
12            output = model(words)
13
14            #output = [sent len, batch size, output dim]
15
16            output = torch.argmax(output, dim=2)
17
18            #output = [sent len, batch size]
19            predict_tags = output.cpu().numpy()
20            true_tags = tags.cpu().numpy()
21
22            true_pred += np.sum(true_tags == predict_tags) & (true_tags != PAD_IDX)
23            num_pred += np.prod(true_tags.shape) - (true_tags == PAD_IDX).sum()
24
25    return true_pred / num_pred

1 lstm_tagger_humor_accuracy = accuracy_model(model, test_iterator)
2 print("LSTM Tagger Humor Accuracy:", lstm_tagger_humor_accuracy * 100, '%')

LSTM Tagger Humor Accuracy: 90.32674118658642 %

```

Вы можете улучшить качество, изменяя параметры модели. Но чтобы добиться нужного качества, вам необходимо взять все выборку, а не только категорию humor.

```

1 brown_tagged_sents = brown.tagged_sents(tagset="universal")

Вам необходимо добиться качества не меньше, чем accuracy = 93 %

Создадим и обучим HMM на всем корпусе
```

```

1 # Приведем слова к нижнему регистру
2 my_brown_tagged_sents = []
3 for sent in brown_tagged_sents:
4     my_brown_tagged_sents.append(list(map(lambda x: (x[0].lower(), x[1]), sent)))
5 my_brown_tagged_sents = np.array(my_brown_tagged_sents)
6
7 from sklearn.model_selection import train_test_split
8 train_sents, test_sents = train_test_split(my_brown_tagged_sents, test_size=0.1, random_state=0)
9
10 hmm = HiddenMarkovModel()
11 hmm.fit(train_sents)
12
13 hmm_tagger_full_accuracy = accuracy_score(hmm, test_sents)
14 print("HMM Tagger Full Accuracy:", hmm_tagger_full_accuracy * 100, '%')

/usr/local/lib/python3.7/dist-packages/ipykernel_launcher.py:5: VisibleDeprecationWarning: Creating an ndarray from ragged nested sequences (which is a list-or-tuple of lists-or-tuples-or ndarrays with different lengths or shapes) is
"""
HMM Tagger Full Accuracy: 96.26295331104619 %

```

Создадим и обучим LSTMTagger на всем корпусе

```

1 pos_data = [list(zip(*sent)) for sent in brown_tagged_sents]
2
3 examples = []
4 for words, tags in pos_data:
5     examples.append(torchtext.legacy.data.Example.fromlist([list(words), list(tags)], fields=[('words', WORD), ('tags', TAG)]))
6
7 dataset = torchtext.legacy.data.Dataset(examples, fields=[('words', WORD), ('tags', TAG)])
8
9 train_data, valid_data, test_data = dataset.split(split_ratio=[0.8, 0.1, 0.1])
10
11 WORD.build_vocab(train_data, min_freq=0.1)
12 TAG.build_vocab(train_data)
13
14 BATCH_SIZE = 512
15
16 train_iterator, valid_iterator, test_iterator = BucketIterator.splits(
17     (train_data, valid_data, test_data),
18     batch_size = BATCH_SIZE,
19     device = device,
20     sort_key=_len_sort_key
21 )

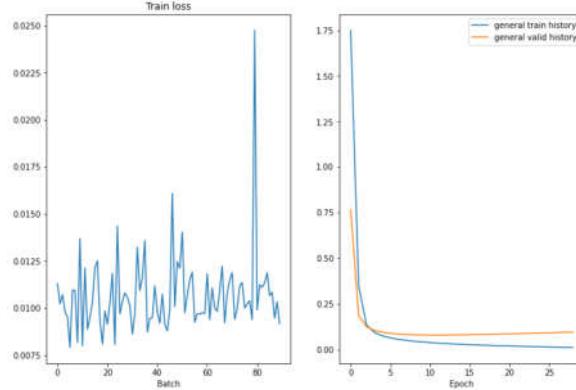
1 INPUT_DIM = len(WORD.vocab)
2 OUTPUT_DIM = len(TAG.vocab)
3 EMB_DIM = 300
4 HID_DIM = 50
5 DROPOUT = 0.5
6 BIDIRECTIONAL = True
7
8 model = LSTMTagger(INPUT_DIM, EMB_DIM, HID_DIM, OUTPUT_DIM, DROPOUT, BIDIRECTIONAL).to(device)
9 model.apply(init_weights)
10
11 PAD_IDX = TAG.vocab.stoi['<pad>']
12 optimizer = optim.Adam(model.parameters())
13 criterion = nn.CrossEntropyLoss(ignore_index = PAD_IDX)
14
15 train_history = []
16 valid_history = []
17
18 N_EPOCHS = 30
19 CLIP = 2
20
21 best_valid_loss = float('inf')
22
23 for epoch in range(N_EPOCHS):
24     start_time = time.time()

```

```

4#     start_time = time.time()
5
6     train_loss = train(model, train_iterator, optimizer, criterion, CLIP, train_history, valid_history)
7     valid_loss = evaluate(model, valid_iterator, criterion)
8
9     end_time = time.time()
10
11    epoch_mins, epoch_secs = epoch_time(start_time, end_time)
12
13    if valid_loss < best_valid_loss:
14        best_valid_loss = valid_loss
15        torch.save(model.state_dict(), 'best-val-model.pt')
16
17    train_history.append(train_loss)
18    valid_history.append(valid_loss)
19    print(f'Epoch: {epoch+1:02} | Time: {epoch_mins}m {epoch_secs}s')
20    print(f'\tTrain loss: {train_loss:.3f} | Train PPL: {math.exp(train_loss):7.3f}')
21    print(f'\tVal. loss: {valid_loss:.3f} | Val. PPL: {math.exp(valid_loss):7.3f}')
22
23

```



```

Epoch: 30 | Time: 0m 5s
  Train Loss: 0.011 | Train PPL:  1.011
  Val. Loss: 0.096 | Val. PPL:  1.101

```

```

1 best_model = LSTMTagger(INPUT_DIM, EMB_DIM, HID_DIM, OUTPUT_DIM, DROPOUT, BIDIRECTIONAL).to(device)
2 best_model.load_state_dict(torch.load('best-val-model.pt'))
3 lstm_tagger_full_accuracy = accuracy_model(best_model, test_iterator)
4 print("LSTM Tagger Full Accuracy:", lstm_tagger_full_accuracy * 100, '%')
5 assert lstm_tagger_full_accuracy * 100 >= 93

```

```
LSTM Tagger Full Accuracy: 97.62193210603455 %
```

Пример решения нашей задачи:

```

1 def print_tags(model, data):
2     model.eval()
3
4     with torch.no_grad():
5         words, _ = data
6         example = torch.LongTensor([WORD.vocab.stoi[i] for i in words]).unsqueeze(1).to(device)
7
8         output = model(example).argmax(dim=-1).cpu().numpy()
9         tags = [TAG.vocab.itos[int(i)] for i in output]
10
11    for token, tag in zip(words, tags):
12        print(f'{token}: {tag}')

1 print_tags(model, pos_data[-1])

```

From	NOUN
what	DET
I	NOUN
was	VERB
able	ADJ
to	ADP
gauge	NOUN
in	ADP
a	DET
swift	ADJ
,	.
greedy	ADJ
glance	NOUN
,	.
the	DET
figure	NOUN
inside	ADP
the	DET
coral-colored	ADJ
boucle	NOUN
dress	NOUN
was	VERB
stupefying	VERB
.	.

#### Справните результаты моделей HiddenMarkov, LstmTagger:

- при обучение на маленькой части корпуса, например, на категории humor
- при обучении на всем корпусе

```

1 print("HMM Tagger Humor Accuracy:", hmm_tagger_humor_accuracy * 100, '%')
2 print("Default Tagger Humor Accuracy:", default_tagger_humor_accuracy * 100, '%')
3 print("NLTK Tagger Humor Accuracy:", nltk_tagger_humor_accuracy * 100, '%')
4 print("RNN Morph Tagger Accuracy:", rnn_morph_tagger_humor_accuracy * 100, '%')
5 print("LSTM Tagger Humor Accuracy:", lstm_tagger_humor_accuracy * 100, '%')
6
7 print("HMM Tagger Full Accuracy:", hmm_tagger_full_accuracy * 100, '%')
8 print("LSTM Tagger Full Accuracy:", lstm_tagger_full_accuracy * 100, '%')

HMM Tagger Humor Accuracy: 88.82847256549678 %
Default Tagger Humor Accuracy: 20.27499764211566 %
NLTK Tagger Humor Accuracy: 89.23302486406328 %
RNN Morph Tagger Accuracy: 62.8274483934758376 %
LSTM Tagger Humor Accuracy: 99.32674118658642 %
HMM Tagger Full Accuracy: 96.26295331104619 %
LSTM Tagger Full Accuracy: 97.62193210603455 %

```

#### Выводы

- LSTM Tagger лучше на малой части корпуса (категория humor) чем HMM Tagger
- LSTM Tagger лучше на всем корпусе чем HMM Tagger

