



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 (HMM).

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

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 руморphy2)
- 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 [Вот тут](#), а про теги — [Вот тут](#)

```
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',
```



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'),  
 ('produced', '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 частотой 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

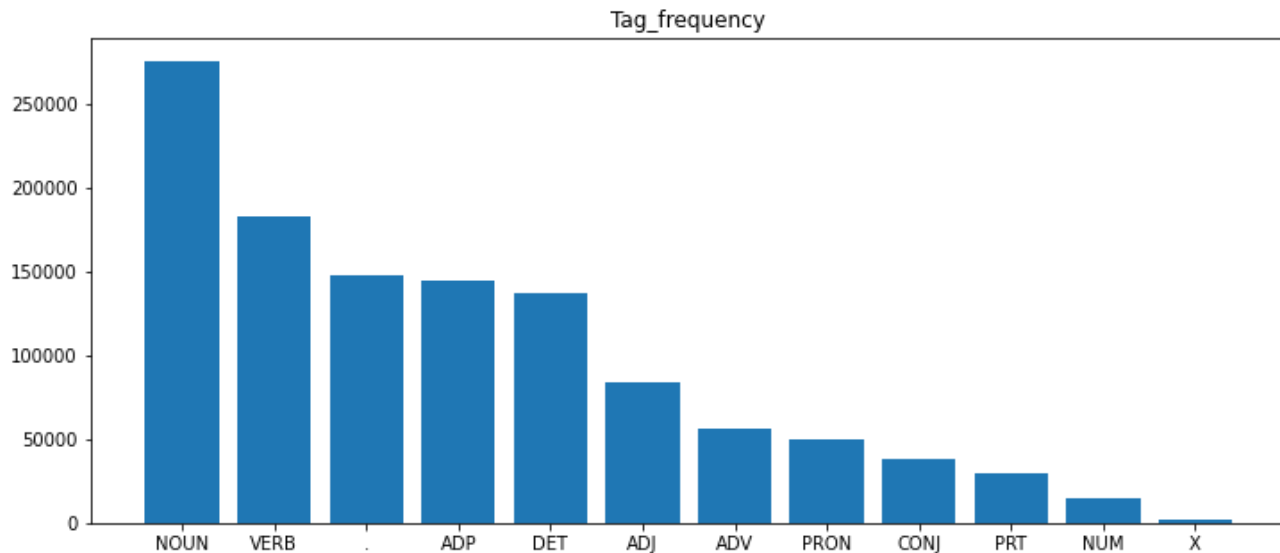
```
5 tag_num = pd.Series(nltk.FreqDist(tags)).sort_values(ascending=False) # тег - кол-во те  
6 word_num = pd.Series(nltk.FreqDist(words)).sort_values(ascending=False) # слово - кол-в
```

Кол-во предложений: 57340

1 tag_num

```
NOUN    275558  
VERB    182750  
.  
147565  
ADP     144766  
DET     137019  
ADJ     83721  
ADV     56239  
PRON    49334  
CONJ    38151  
PRT     29829  
NUM     14874  
X        1386  
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()
```



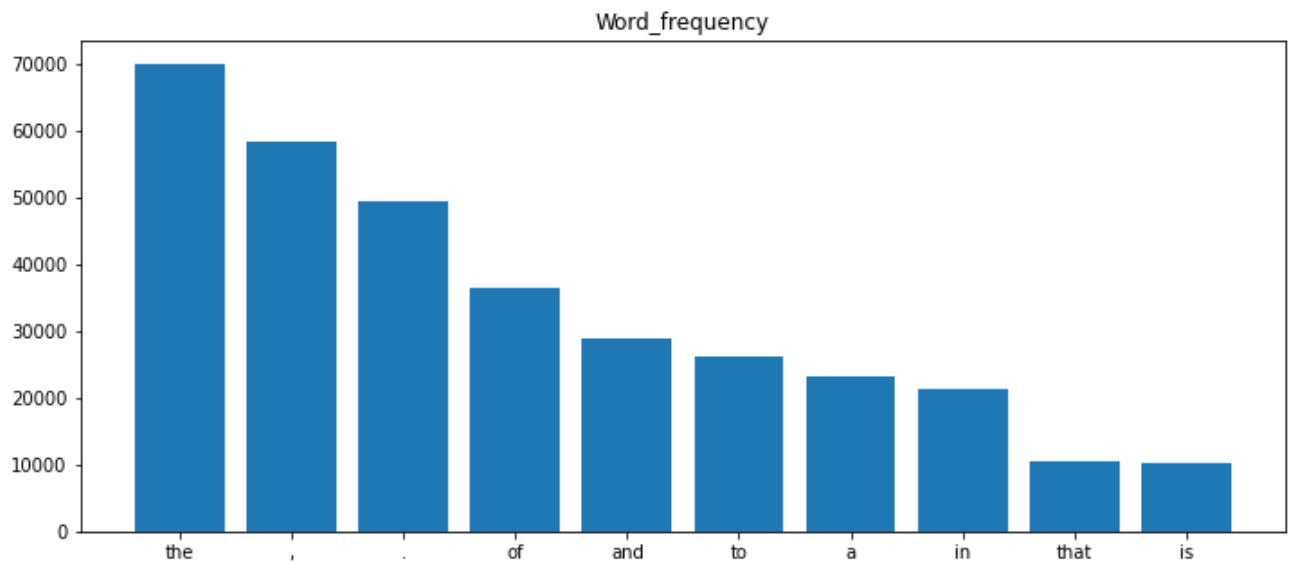
1 word_num[:5]

```
the      69971  
,        58334  
.  
49346  
of       36412  
and      28853  
dtype: int64
```

```

1 plt.figure(figsize=(12, 5))
2 plt.bar(word_num.index[:10], word_num.values[:10])
3 plt.title("Word_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 words_with_most_popular_tag = [word for (word, tag) in brown_tagged_words if tag == tag
2 words_with_most_popular_tag_num = pd.Series(nltk.FreqDist(words_with_most_popular_tag))

```

```
1 words_with_most_popular_tag_num.index[0]
```

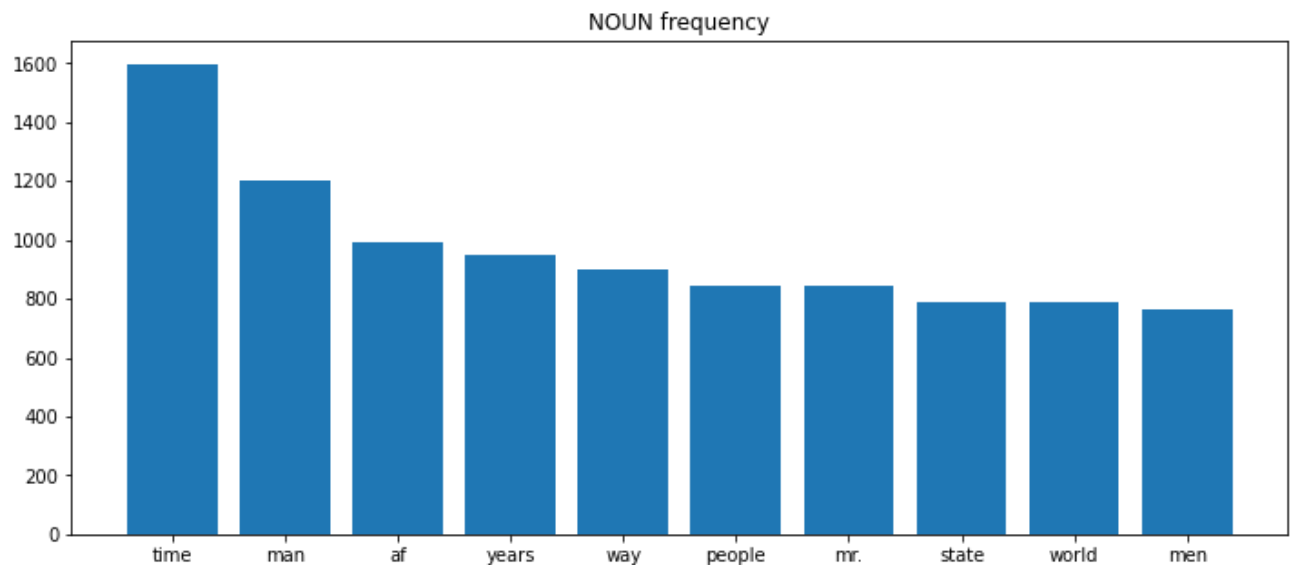
'time'

```

1 plt.figure(figsize=(12, 5))
2 plt.bar(words_with_most_popular_tag_num.index[:10], words_with_most_popular_tag_num.val
3 plt.title(f"{tag_num.index[0]} frequency")
4 plt.show()

```

```
5 print()
6 print(f'Most popular {tag_num.index[0]}:', words_with_most_popular_tag_num.index[0])
```



Most popular NOUN: time

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 brown_tagged_sents = brown.tagged_sents(tagset="universal")
3
```

```

4 # Приведем слова к нижнему регистру
5 my_brown_tagged_sents = []
6 for sent in brown_tagged_sents:
7     my_brown_tagged_sents.append(list(map(lambda x: (x[0].lower(), x[1]), sent)))
8 my_brown_tagged_sents = np.array(my_brown_tagged_sents)
9
10 from sklearn.model_selection import train_test_split
11 train_sents, test_sents = train_test_split(my_brown_tagged_sents, test_size=0.1, random

```

/usr/local/lib/python3.7/dist-packages/ipykernel_launcher.py:8: VisibleDeprecationWar

```
1 len(train_sents)
```

51606

```
1 len(test_sents)
```

5734

▼ 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
- $a_{i,j} = p(s_j | s_i)$ - вероятность того, что, находясь в скрытом состоянии s_i , мы попадем в состояние s_j (элемент матрицы A)
- $b_{k,j} = p(o_k | s_j)$ - 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 размеченным данным (последовательности пар слово+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-тег - 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.

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

$Q_{t,s}$ можно восстановить по `argmax`-ам.

```

1 class HiddenMarkovModel:
2     def __init__(self):
3
4         pass
5
6     def fit(self, train_tokens_tags_list):
7         """
8         train_tokens_tags_list: массив предложений пар слово-тег (выборка для train)
9         """
10        tags = [tag for sent in train_tokens_tags_list
11                for (word, tag) in sent]
12        words = [word for sent in train_tokens_tags_list for (word, tag) in sent]
13
14        tagged_words = [(word, tag) for sent in train_tokens_tags_list
15                        for (word, tag) in sent]
16
17        tag_num = pd.Series(nltk.FreqDist(tags)).sort_values(ascending=False)
18        word_num = pd.Series(nltk.FreqDist(words)).sort_values(ascending=False)
19
20        self.tags = tag_num.index
21        self.words = word_num.index
22
23        mpv = self.words[0] # most popular word
24        mpt = self.tags[0]  # most popular tag
25        # words with most popular tag
26        words_with_most_popular_tag = [word for (word, tag) in tagged_words if tag == mpt]

```



```

27     # most popular word
28     mpwmpmt = pd.Series(nltk.FreqDist(words_with_most_popular_tag)).sort_values(asc=
29
30
31     self.most_popular_tag = mpt
32     self.most_popular_word = mpv
33     self.most_popular_word_from_most_popular_tag = mpwmpmt
34
35
36
37
38
39     A = pd.DataFrame({'{}'.format(tag) : [0] * len(tag_num) for tag in tag_num.inde
40     B = pd.DataFrame({'{}'.format(tag) : [0] * len(word_num) for tag in tag_num.ind
41
42     # Вычисляем матрицу A и B по частотам слов и тегов
43
44     # sent - предложение
45     # sent[i][0] - i слово в этом предложении, sent[i][1] - i тег в этом предложени
46     for sent in train_tokens_tags_list:
47         for i in range(len(sent)):
48             B.loc[sent[i][0], sent[i][1]] += 1 # текущая i-пара слово-тег (обновите
49             if len(sent) - 1 != i: # для последнего тега нет следующего тега
50                 A.loc[sent[i][1], sent[i + 1][1]] += 1 # пара тег-тег
51
52
53     # переходим к вероятностям
54
55     # нормируем по строке, то есть по всем всевозможным следующим тегам
56     A = A.divide(A.sum(axis=1), axis=0)
57
58     # нормируем по столбцу, то есть по всем всевозможным текущим словам
59     B = B / np.sum(B, axis=0)
60
61     self.A = A
62     self.B = B
63
64     return self
65
66
67 def predict(self, test_tokens_list):
68     """
69     test_tokens_list : массив предложений пар слово-тег (выборка для test)
70     """
71     predict_tags = OrderedDict({i : np.array([]) for i in range(len(test_tokens_lis
72
73     for i_sent in range(len(test_tokens_list)):
74
75         current_sent = [] # текущее предложение
76         len_sent = len(test_tokens_list[i_sent]) # длина предложения
77
78         q = np.zeros(shape=(len_sent + 1, len(self.tags)))
79         q[0] = 1 # нулевое состояние (равномерная инициализация по всем s)
80         back_point = np.zeros(shape=(len_sent + 1, len(self.tags))) # # argmax
81

```

```

82         for t in range(len_sent):
83
84             # если мы не встречали такое слово в обучении, то вместо него будет
85             # самое популярное слово с самым популярным тегом (вопрос 2)
86             current_word = test_tokens_list[i_sent][t]
87             if current_word not in self.words:
88                 current_sent.append(self.most_popular_word_from_most_popular_tag)
89             else:
90                 current_sent.append(current_word)
91
92             # через max выбираем следующий тег
93             for i_s in range(len(self.tags)):
94
95                 s = self.tags[i_s]
96
97                 # формула (1)
98                 q[t + 1][i_s] = np.max(q[t] *
99                     self.A.loc[:, s] *
100                     self.B.loc[current_sent[t], s])
101
102
103             # argmax формула(1)
104
105             # argmax, чтобы восстановить последовательность тегов
106             back_point[t + 1][i_s] = (q[t] * self.A.loc[:, s] *
107                 self.B.loc[current_sent[t], s]).reset_index()[s].idxmax() # инде
108
109         back_point = back_point.astype('int')
110
111         # выписываем теги, меняя порядок на реальный
112         back_tag = deque()
113         current_tag = np.argmax(q[len_sent])
114         for t in range(len_sent, 0, -1):
115             back_tag.appendleft(self.tags[current_tag])
116             current_tag = back_point[t, current_tag]
117
118         predict_tags[i_sent] = np.array(back_tag)
119
120
121     return predict_tags

```

Обучите скрытую марковскую модель:

```

1 my_model = HiddenMarkovModel()
2 my_model.fit(train_sents)

<__main__.HiddenMarkovModel at 0x7facc6546b50>

```

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'],
2          ['a', 'cat', 'and', 'a', 'dog'],
3          ['I', 'have', 'a', 'television'],
4          ['My', 'favourite', 'character']]
5
6 preds = my_model.predict(sents)
7 preds

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 my_model.predict(['can'])

OrderedDict([(0, array(['VERB'], dtype='<U4'))])

```

Удивительно, если написать Can, то модель считает его существительным

▼ Вопрос 4:

- Какой тег вы получили для слова favourite?

```

1 my_model.predict(['favourite'])

OrderedDict([(0, array(['NOUN'], dtype='<U4'))])

```

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     i = 0
5
6     for sent in sents:
7         tags = np.array([tag for (word, tag) in sent])
8         words = np.array([word for (word, tag) in sent])
9         preds = np.array(model.predict(words)[0])

```

```

10
11
12     '''your code'''
13     match = (preds == tags)
14
15
16     true_pred += match.sum()
17     num_pred += len(match)
18     print("Accuracy:", true_pred / num_pred * 100, '%')
19     return (true_pred / num_pred)

```

```

1 accu = accuracy_score(my_model, test_sents)

```

Accuracy: 96.26470820500671 %

▼ Вопрос 5:

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

```

1 np.round(accu, 1)

```

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('NOUN')

```

```

1 def default_accuracy_score(model, sents):
2     true_pred = 0
3     num_pred = 0
4
5     for sent in sents:
6         tags = np.array([tag for (word, tag) in sent])
7
8         words = [word for (word, tag) in sent]
9         words_and_preds = model.tag(words)
10        preds = np.array([pred for (word, pred) in words_and_preds])

```



```
Collecting dawg-python>=0.7.1
  Downloading DAWG_Python-0.7.2-py2.py3-none-any.whl (11 kB)
Requirement already satisfied: joblib>=0.11 in /usr/local/lib/python3.7/dist-packages
Requirement already satisfied: typing-extensions>=3.6.4 in /usr/local/lib/python3.7/dist-packages
Requirement already satisfied: zipp>=0.5 in /usr/local/lib/python3.7/dist-packages
Building wheels for collected packages: rnnmorph, russian-tagsets
  Building wheel for rnnmorph (setup.py) ... done
  Created wheel for rnnmorph: filename=rnnmorph-0.4.1-py3-none-any.whl size=19746379
  Stored in directory: /root/.cache/pip/wheels/b9/e8/8b/b4639d66ee5373f0db8dbe9a2ea94
  Building wheel for russian-tagsets (setup.py) ... done
  Created wheel for russian-tagsets: filename=russian_tagsets-0.6-py3-none-any.whl si
  Stored in directory: /root/.cache/pip/wheels/5e/b4/26/9c17a7cdcfc6b8cf43111312f3e7c
Successfully built rnnmorph russian-tagsets
Installing collected packages: pymorphy2-dicts-ru, dawg-python, russian-tagsets, pymc
Successfully installed dawg-python-0.7.2 jsonpickle-2.0.0 pymorphy2-0.9.1 pymorphy2-c
```

```
1 from nltk.tag.mapping import map_tag
```

```
1 import nltk
2 nltk.download('averaged_perceptron_tagger')
3 # nltk.pos_tag(..., tagset='universal')
4
```

```
[nltk_data] Downloading package averaged_perceptron_tagger to
[nltk_data] /root/nltk_data...
[nltk_data] Unzipping taggers/averaged_perceptron_tagger.zip.
True
```

```
1 from rnnmorph.predictor import RNNMorphPredictor
2 predictor = RNNMorphPredictor(language="en")
```

```
[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-
[nltk_data] 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 kernels since it doesn't m
WARNING:tensorflow:Layer LSTM_1_backward will not use cuDNN kernels since it doesn't
WARNING:tensorflow:Layer LSTM_0 will not use cuDNN kernels since it doesn't meet the
WARNING:tensorflow:Layer LSTM_0 will not use cuDNN kernels since it doesn't meet the
WARNING:tensorflow:Layer LSTM_0 will not use cuDNN kernels since it doesn't meet the
```

```
1 # predictor.predict(np.array(['Hello' 'sadaf']))
2 predictor.predict(['I', 'bear', 'a', 'bag'])[2].pos

'DET'
```

▼ Бонус 7:

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

```

1 def nltk_accuracy_score(tagset, sents):
2     true_pred = 0
3     num_pred = 0
4
5     for sent in sents:
6         tags = np.array([tag for (word, tag) in sent])
7
8         words = [word for (word, tag) in sent]
9         words_and_preds = nltk.pos_tag(words, tagset=tagset)
10        preds = np.array([pred for (word, pred) in words_and_preds])
11        # print(preds)
12
13        match = (tags == preds)
14        #print(match)
15
16        true_pred += match.sum()
17        num_pred += len(preds)
18
19    acc = true_pred / num_pred
20    print("Accuracy:", np.round(acc * 100, 1), '%')
21    return acc

```

```

1 nltk_acc_score = nltk_accuracy_score('universal', test_sents)
2 nltk_acc_score

```

```

Accuracy: 89.2 %
0.8922392486406328

```

```

1 def accuracy_score_r(model, sents):
2     true_pred = 0
3     num_pred = 0
4
5     for sent in sents:
6         tags = np.array([tag for (word, tag) in sent])
7         words = [word for (word, tag) in sent]
8
9         full_preds = predictor.predict(words)
10
11        preds = np.array([full_preds[i].pos for i in range(len(full_preds))])
12
13
14        # preds = np.array(list(model.predict(word)[0].pos for word in words))
15
16        # print(zip(preds[:10], tags[:10]))
17        # for (k, m) in zip(preds[:10], tags[:10]):

```

```

18         #     print(k, m)
19
20         '''your code'''
21         match = (preds == tags)
22
23
24         true_pred += match.sum()
25         num_pred += len(match)
26         print("Accuracy:", true_pred / num_pred * 100, '%')
27         return (true_pred / num_pred)

```

```
1 accuracy_score_r(predictor, test_sents)
```

```

Accuracy: 62.827483934750376 %
0.6282748393475037

```

▼ BiLSTMTagger

▼ Подготовка данных

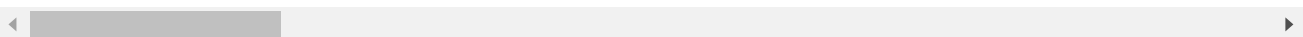
Изменим структуру данных

```

1 pos_data = [list(zip(*sent)) for sent in brown_tagged_sents]
2 print(pos_data[0])

```

```
[('The', 'Fulton', 'County', 'Grand', 'Jury', 'said', 'Friday', 'an', 'investigation
```



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

```

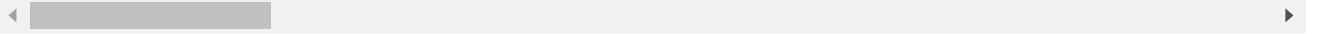
1 from torchtext.legacy.data import Field, BucketIterator
2 import torchtext
3
4 # наши поля
5 WORD = Field(lower=True)
6 TAG = Field(unk_token=None) # все токены нам известны
7
8 # создаем примеры
9 examples = []
10 for words, tags in pos_data:
11     examples.append(torchtext.legacy.data.Example.fromlist([list(words),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

```

Вот один наш пример:

```
1 print(vars(examples[0]))
```

```
{'words': ['the', 'fulton', 'county', 'grand', 'jury', 'said', 'friday', 'an', 'inves
```

Теперь формируем наш датасет

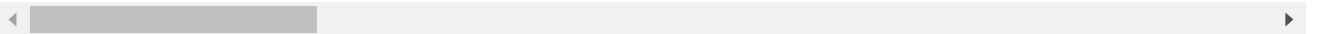
```
1 # кладем примеры в наш датасет
2 dataset = torchtext.legacy.data.Dataset(examples, fields=[('words', WORD), ('tags', TAG)
3
4 train_data, valid_data, test_data = dataset.split(split_ratio=[0.8, 0.1, 0.1])
5
6 print(f"Number of training examples: {len(train_data.examples)}")
7 print(f"Number of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

```
Number of training examples: 45872
Number of validation examples: 5734
Number of testing examples: 5734
```

Построим словари. Параметр `min_freq` выберете сами. При построении словаря используем только **train**

```
1 WORD.build_vocab(train_data, min_freq=10)
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: 7325
Unique tokens in target (en) vocabulary: 13
['<unk>', 'away', 'close', 'bring', 'southern', 'leaders', 'units', 'robert', 'signs
['<pad>', 'NOUN', 'VERB', '.', 'ADP', 'DET', 'ADJ', 'ADV', 'PRON', 'CONJ', 'PRT', 'NL
```



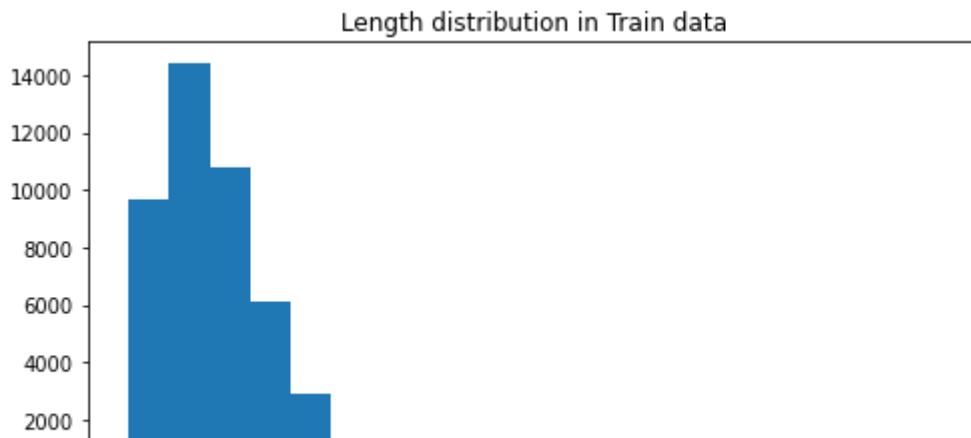
```
1 print(vars(train_data.examples[9]))

{'words': ['the', "women's", 'faces', 'had', 'hardened', 'after', 'my', 'statement',
```



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);
```



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 = 128
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 )

1 # посмотрим на количество батчей
2 list(map(len, [train_iterator, valid_iterator, test_iterator]))

[359, 45, 45]
```

▼ 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=
4         super().__init__()
5
6
7         self.embeddings = nn.Embedding(num_embeddings=input_dim, embedding_dim=emb_dim)
8         self.dropout = nn.Dropout(p=dropout)
9
10        self.rnn = nn.LSTM(input_size=emb_dim, hidden_size=hid_dim,
11                            dropout=dropout, bidirectional=bidirectional)
12        # если bidirectional, то предсказываем на основе конкатенации двух hidden
13        self.tag = nn.Linear((1 + bidirectional) * hid_dim, output_dim)
14
15    def forward(self, sent):
16
17        #sent = [sent len, batch size]
18
19        # не забываем применить dropout к embedding
20        embedded =self.embeddings(sent)
21
22        output, _ = self.rnn(self.dropout(embedded))
23        #output = [sent len, batch size, hid dim * n directions]
24
25        prediction = self.tag(output)
26
27        return prediction
28
29 # параметры модели
30 INPUT_DIM = len(WORD.vocab) # количество эмбеддингов
31 OUTPUT_DIM = len(TAG.vocab) # количество классифицируемых тегов
32 EMB_DIM = 300
33 HID_DIM = 16
34 DROPOUT = 0.2
35 BIDIRECTIONAL =True
36
37 model = LSTMTagger(input_dim=INPUT_DIM, emb_dim=EMB_DIM,
38                    hid_dim=HID_DIM, output_dim=OUTPUT_DIM,
39                    dropout=DROPOUT, bidirectional=BIDIRECTIONAL).to(device)
40
41 # инициализируем веса
42 def init_weights(m):
43     for name, param in m.named_parameters():
44         nn.init.uniform_(param, -0.08, 0.08)
45
46 model.apply(init_weights)

```

```

/usr/local/lib/python3.7/dist-packages/torch/nn/modules/rnn.py:65: UserWarning: dropout
"num_layers={}".format(dropout, num_layers))

```

```

LSTMTagger(
  (embeddings): Embedding(7325, 300)
  (dropout): Dropout(p=0.2, inplace=False)
  (rnn): LSTM(300, 16, dropout=0.2, bidirectional=True)
  (tag): Linear(in_features=32, out_features=13, bias=True)
)

```

Подсчитаем количество обучаемых параметров нашей модели

```
1 def count_parameters(model):
2     return sum(p.numel() for p in model.parameters())
3
4 print(f'The model has {count_parameters(model):,} trainable parameters')
```

The model has 2,238,633 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        '''your code'''
12        words = batch.words.to(device)
13        tags = batch.tags.to(device)
14
15        optimizer.zero_grad()
16
17        output = model(words)
18
19        #tags = [sent len, batch size]
20        #output = [sent len, batch size, output dim]
21
22        output = output.view(-1, output.shape[2])
23        tags = tags.view(-1)
24
25        #tags = [sent len * batch size]
26        #output = [sent len * batch size, output dim]
27
28        loss = criterion(output, tags)
29
30        loss.backward()
31
32        # Gradient clipping(решение проблемы взрыва градиентов), clip - максимальная норма
33        torch.nn.utils.clip_grad_norm_(model.parameters(), max_norm=clip)
34
35        optimizer.step()
36
37        epoch_loss += loss.item()
38
39        history.append(loss.cpu().data.numpy())
40        if (i+1)%10==0:
```

```

41         fig, ax = plt.subplots(nrows=1, ncols=2, figsize=(12, 8))
42
43         clear_output(True)
44         ax[0].plot(history, label='train loss')
45         ax[0].set_xlabel('Batch')
46         ax[0].set_title('Train loss')
47
48         if train_history is not None:
49             ax[1].plot(train_history, label='general train history')
50             ax[1].set_xlabel('Epoch')
51         if valid_history is not None:
52             ax[1].plot(valid_history, label='general valid history')
53         plt.legend()
54
55         plt.show()
56
57
58     return epoch_loss / len(iterator)
59
60 def evaluate(model, iterator, criterion):
61     model.eval()
62
63     epoch_loss = 0
64
65     history = []
66
67     with torch.no_grad():
68
69         for i, batch in enumerate(iterator):
70
71             '''your code'''
72             words = batch.words.to(device)
73             tags = batch.tags.to(device)
74
75             output = model(words)
76
77             #tags = [sent len, batch size]
78             #output = [sent len, batch size, output dim]
79
80             output = output.view(-1, output.shape[2])
81             tags = tags.view(-1)
82
83             #tags = [sent len * batch size]
84             #output = [sent len * batch size, output dim]
85
86             loss = criterion(output, tags)
87
88             epoch_loss += loss.item()
89
90     return epoch_loss / len(iterator)
91
92 def epoch_time(start_time, end_time):
93     elapsed_time = end_time - start_time
94     elapsed_mins = int(elapsed_time / 60)

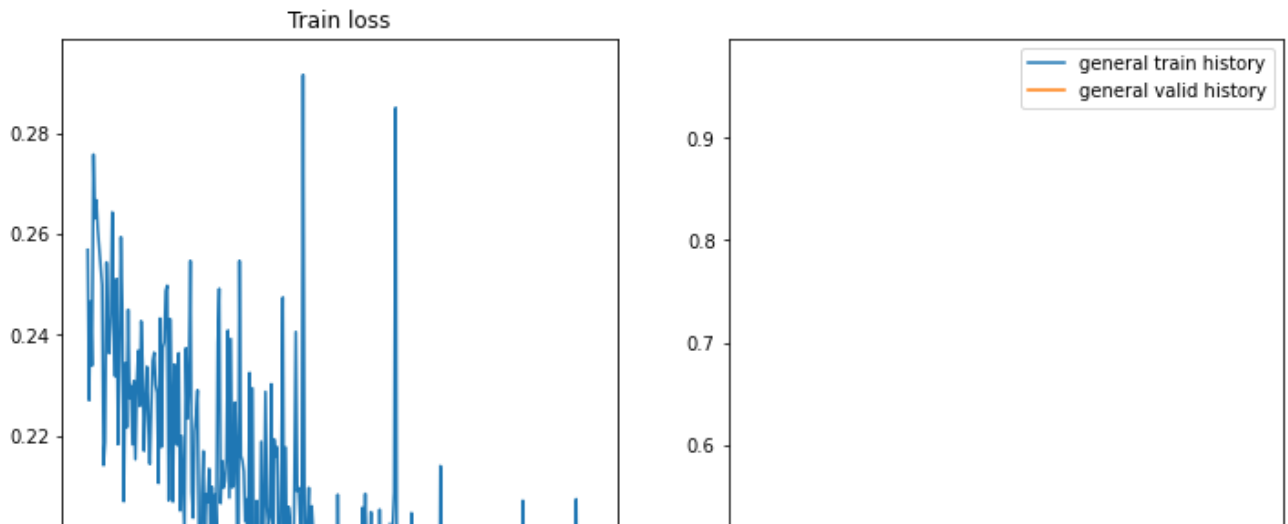
```

```

95     elapsed_secs = int(elapsed_time - (elapsed_mins * 60))
96     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 = 2
13 CLIP = 1
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)
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}m {epoch_secs}s')
35     print(f'\tTrain Loss: {train_loss:.3f} | Train PPL: {math.exp(train_loss):7.3f}')
36     print(f'\tVal. Loss: {valid_loss:.3f} | Val. PPL: {math.exp(valid_loss):7.3f}')

```



▼ 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             '''your code'''
10            words = batch.words.to(device)
11            tags = batch.tags.to(device)
12
13            output = model(words)
14
15            #output = [sent len, batch size, output dim]
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 round(true_pred / num_pred * 100, 3)

```

```
1 print("Accuracy:", accuracy_model(model, test_iterator), '%')
```

Accuracy: 95.707 %

```
1 print("Accuracy:", accuracy_model(model, test_iterator), '%')
```

Accuracy: 95.707 %

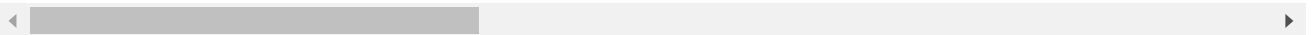
Вы можете улучшить качество, изменяя параметры модели. Но чтобы добиться нужного качества, вам необходимо взять все выборку, а не только категорию `humor`.

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

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

```
1 best_model = LSTMTagger(INPUT_DIM, EMB_DIM, HID_DIM, OUTPUT_DIM, DROPOUT, BIDIRECTIONAL
2 best_model.load_state_dict(torch.load('best-val-model.pt'))
3 assert accuracy_model(best_model, test_iterator) >= 93
```

```
/usr/local/lib/python3.7/dist-packages/torch/nn/modules/rnn.py:65: UserWarning: dropout
  "num_layers={}".format(dropout, num_layers))
```



Пример решение 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[elem] for elem in words]).unsqueeze
7
8         output = model(example).argmax(dim=-1).cpu().numpy()
9         tags = [TAG.vocab.itos[int(elem)] for elem in output]
10
11     for token, tag in zip(words, tags):
12         print(f'{token:15s}{tag}')
```

```
1 print_tags(model, pos_data[-1])
```

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

boucle	NOUN
dress	NOUN
was	VERB
stupefying	ADJ
.	.

Сравните результаты моделей HiddenMarkov, LstmTagger:

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

▼ HiddenMarkov

На категории humor 88 %

На всем корпусе 96 %(но считалось целый час)

LstmTagger

На категории humor - 86 %

На всем корпусе - 95 %

1