

2. ВТОРОЙ ЭТАП

Задачи второго этапа

Задача 3.0.1. Белые или чёрные? (15 баллов)

Winning and losing isn't everything;

sometimes, the journey is just as important as the outcome.

Alex Morgan

В первой задаче вам предстоит анализировать шахматные партии. Вам дан датасет, в котором содержатся описания шахматных партий за некоторый период времени с одного сайта для онлайн-игры в шахматах. Вам гарантируется, что каждая партия закончилась победой либо чёрных, либо белых. Ваша задача будет состоять в том, чтобы для каждой партии ответить, какая сторона победила.

Формат входных данных

Вам надо будет скачать датасет с играми (https://drive.google.com/file/d/1IsebLv8-AcLoBbIb13gCVJ7cu_oo0eQV/view). Он имеет типичный для машинного обучения формат .csv. Про каждую партию вам известна следующая информация:

1. 'id' – уникальный идентификатор партии
2. 'rated' – True/False – была ли партия рейтинговой
3. 'created_at' – время начала партии
4. 'last_move_at' – во сколько был совершён последний ход
5. 'turns' – сколько ходов было в партии
6. 'white_id' – уникальный идентификатор игрока белыми
7. 'white_rating' – рейтинг игрока белыми
8. 'black_id' – уникальный идентификатор игрока чёрными
9. 'black_rating' – рейтинг игрока чёрными

Формат выходных данных

В систему вам необходимо сдать файл, в котором будет один столбец с заголовком 'winner', в каждой ячейке которого будет написано 'black' или 'white'. Пример ссылки. (<https://drive.google.com/file/d/1NE48oM5IuBGLbrJ1Dk9Noti4mdsммj8X/view?usp=sharing>)

Система оценки

За эту задачу вы можете получить 0 или 15 баллов: 15 – в случае, если ваш ответ полностью верен, 0 – в противном случае.

Pandas - самая популярная библиотека для работы с табличными данными при анализе данных.

Пример программы-решения

Ниже представлено решение на языке Python3

```
1 import pandas as pd
2
3 df = pd.read_csv('input.csv', index_col=0)
4 ans_bool = df.turns.apply(lambda x: 'white' if x % 2 else 'black')
5 answer = pd.DataFrame({'winner' : ans_bool})
6 answer.to_csv('answer.csv', index=False)
```

Задача 3.0.2. Это легковушка? Это автобус? (15 баллов)

В этой задаче вам предстоит по данным GPS-навигации определять тип средства передвижения – это автобус или легковой автомобиль?

Формат входных данных

В качестве тренировочной выборки вам даны некоторые данные GPS-навигации о движении транспорта в Красноярске.

<https://docs.google.com/spreadsheets/d/1QdZJT05NiqgCDPYrPHJBWCK2Se0RK2Fo6HN-ODT7cgE/edit?usp=sharing>

В качестве признаков имеются средняя скорость, время в пути, пройденная дистанция, rating – оценочный параметр, выражающий насколько данное средство передвижения удобно для перемещения по городу (3 – удобно, 2 – нормально, 1 – плохо), а также ответ – 'car' или 'bus'.

Помимо этого вам дана тестовая выборка (test.csv), набор столбцов которой отличается от обучающей только отсутствием столбца car_or_bus. Как вы уже могли догадаться, ваша задача – обучиться на обучающей выборке и научиться предсказывать для тестовой выборки тип каждого из указанных в ней средств передвижения (car или bus).

Формат выходных данных

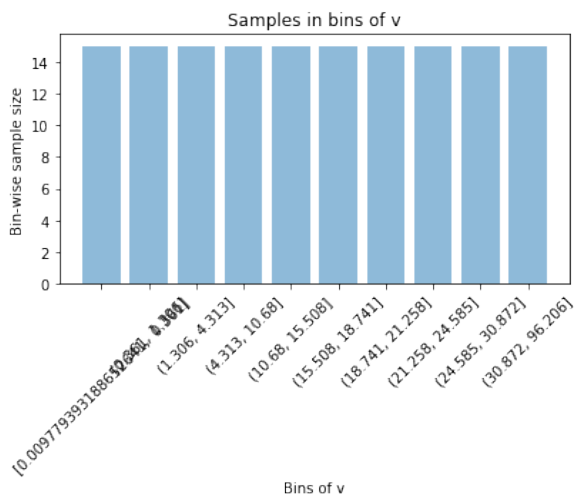
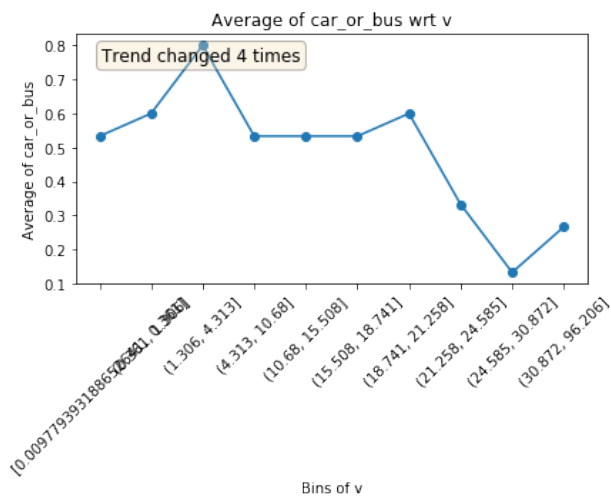
Вам нужно сдать в систему файл, в котором для каждой строки из тестовой выборки указан предполагаемый вами ответ (car/bus). Пример ссылки. (https://drive.google.com/file/d/1dM0Jiq0_UU7ey4o2dcVjhS23pfbKaea5/view?usp=sharing)

Система оценки

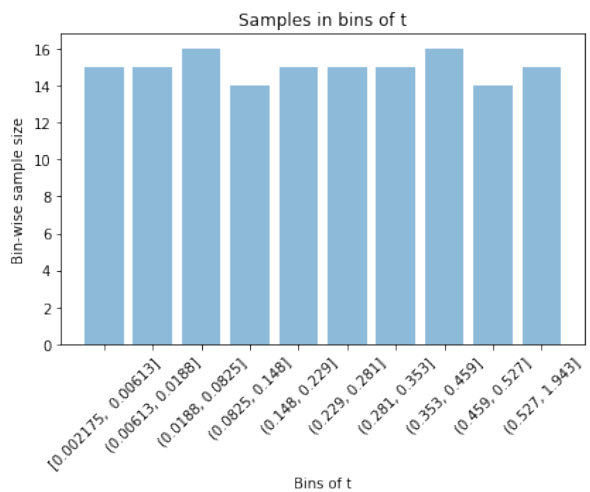
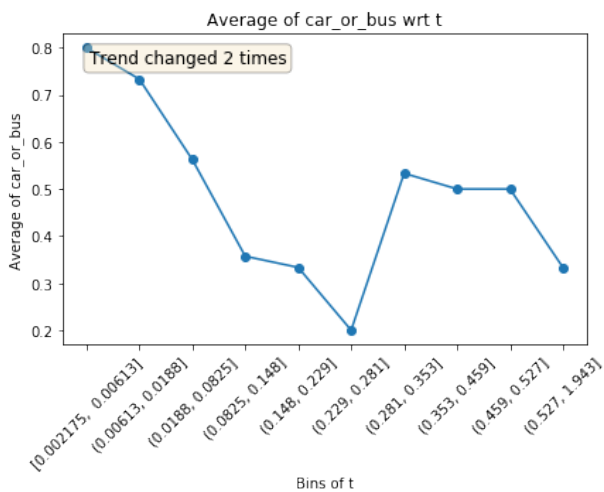
Чем больше правильных ответов вы дадите, тем больше баллов вы заработаете. Назовём accuracy_score величину, равную отношению данных правильно ответов к общему количеству строк в test.csv.

Тогда:

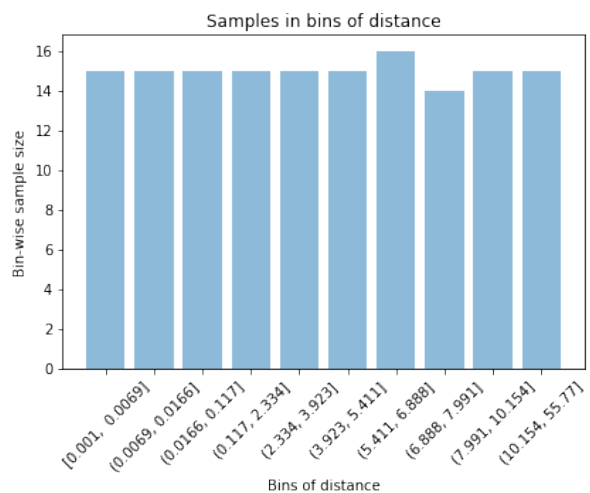
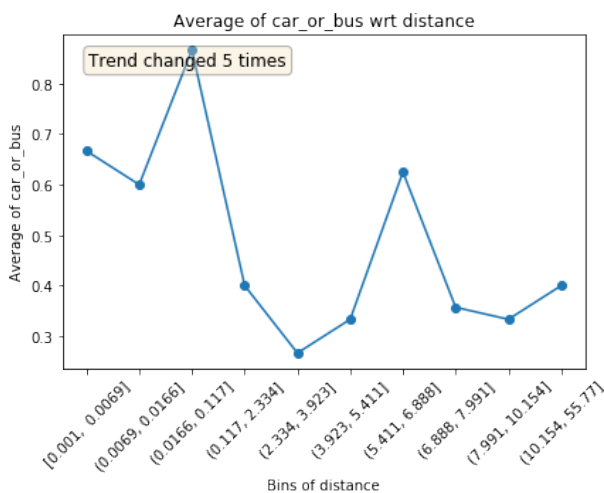
Plots for v



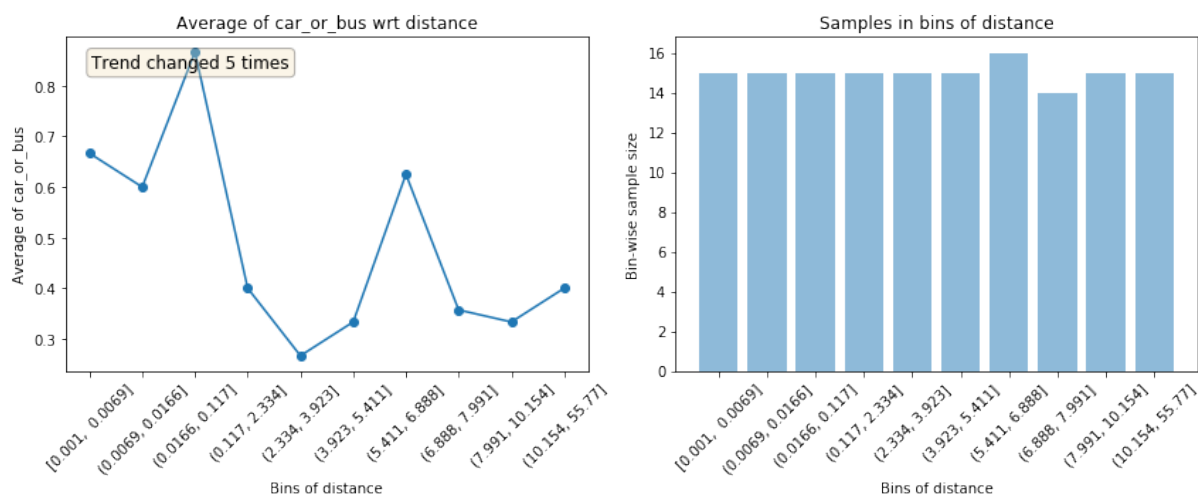
Plots for t



Plots for distance



Plots for rating



Все признаки имеют зависимость с классом, поэтому не стоит их выбрасывать.

Признак distance это $t \cdot v$ (в чем легко убедиться, перемножив столбцы). Он не вносит новых данных и можно попробовать его убрать. В результате эксперимента мы выяснили, что без признака distance точность всех классификаторов становится ниже, а значит, он важен для представления данных.

```

1 # Разделим данные на признаки и метки класса
2 X = data.iloc[:, :-1].values
3 y = data.iloc[:, -1].values
4
5 # Мы будем использовать модели, для которых важно, чтобы признаки были численно
6 # одного порядка, поэтому мы нормализуем их (вычтем среднее и поделим на
7 # стандартное отклонение).
8 # Потом normalizer мы будем использовать для тестовых данных
9 normalizer = StandardScaler()
10 X = normalizer.fit_transform(X)

```

Так как у нас очень мало данных, то мы будем использовать особую стратегию для оценки моделей - Leave One Out. Мы будем выделять одно наблюдение из всего датасета, обучать модель на оставшихся наблюдениях, а затем этой моделью предсказывать метку для выделенного наблюдения.

Такая стратегия нужна для того, чтобы при валидации модель предсказывала метку только для наблюдений, которые она не видела в процессе тренировки, и при этом использовала максимум данных для обучения. В решении следующей задачи мы будем разделять выборку на две части - обучающую и валидационную, не используя такую сложную стратегию, потому что данных намного больше.

```

1 def loo_accuracy(model):
2     """
3     Производит оценку точности модели на датасете X, y (внешние переменные),
4     используя стратегию Leave One Out. Возвращает сре
5     """
6     from sklearn.model_selection import LeaveOneOut
7
8     loo = LeaveOneOut()

```

```

9     correct = 0
10    for train_ind, val_ind in loo.split(X, y):
11        X_train, y_train = X[train_ind], y[train_ind]
12        X_val, y_val = X[val_ind], y[val_ind]
13
14        model.fit(X_train, y_train)
15        correct += (model.predict(X_val) == y_val)[0]
16
17    return correct / len(X)

```

SVM

```

1  svm = SVC(gamma='auto')
2
3  # Выведем точность модели
4  loo_accuracy(svm)

```

Логистическая регрессия

```

1  # Инициализируем модель и обучим на данных
2  log_reg = LogisticRegression(C=0.1, solver='lbfgs')
3
4  # Выведем точность модели
5  loo_accuracy(log_reg)

```

N ближайших соседей

```

1  # Найдем лучшее значение для n_neighbors
2  best_n_neighbors = 0
3  best_acc = 0
4
5  for n_neighbors in range(1, 20):
6      # Инициализируем модель и обучим на данных
7      model = KNeighborsClassifier(n_neighbors=n_neighbors)
8
9      # Найдем точность модели
10     acc = loo_accuracy(model)
11
12     if acc > best_acc:
13         best_n_neighbors = n_neighbors
14         best_acc = acc
15
16     print('Best accuracy: {:.2f} for {} neighbors'.format(float(best_acc),
17         best_n_neighbors))

```

Градиентный бустинг

```

1  # Найдем лучшее значение для n_estimators
2  best_n_estimators = 0
3  best_acc = 0
4
5  for n_estimators in range(1, 30):
6      # Инициализируем модель и обучим на данных
7      model = GradientBoostingClassifier(n_estimators=n_estimators,
8          learning_rate=0.1)
9

```

```

10     # Найдем точность на валидации
11     acc = loo_accuracy(model)
12
13     if acc > best_acc:
14         best_n_neighbors = n_neighbors
15         best_acc = acc
16
17 print('Best accuracy: {:.2f} for {} estimators'.format(float(best_acc),
18         best_n_neighbors))

```

Наивный байес

```

1  # Инициализируем модель и обучим на данных
2  model = GaussianNB()
3
4  # Выведем точность модели
5  loo_accuracy(model)

```

Итоговая модель

Так как N ближайших соседей оказался лучшим, мы будем использовать именно его с лучшим `n_neighbors`

```

1  # Инициализируем модель, используя best_n_neighbors = 4
2  model = KNeighborsClassifier(n_neighbors=best_n_neighbors)
3
4  # обучим на всем датасете
5  model.fit(X, y)

```

Получение ответа

```

1  # Загрузим тестовые данные
2  data_test = pd.read_csv('DATASET2.csv')
3
4  # Нормализуем их, используя тот же normalizer, что и для обучающей выборки
5  X_test = normalizer.transform(data_test.values)
6
7  # Запишем результат в pandas.Series
8  result = pd.Series(model.predict(X_test))
9
10 # Преобразуем 0 и 1 в метки классов, сделав словарь с обратным маппингом
11 inv_map = {v: k for k, v in mapping.items()}
12 result = result.map(inv_map)
13 result
14
15 # Запишем результат в файл
16 result.to_csv('second_task_res.csv', index=False)

```

Если вы никогда раньше не работали с `pandas` и `sklearn`, то вы можете начать свое изучение с этой статьи (<https://habr.com/post/202090/>). Так же очень полезно пользоваться документацией по `sklearn` и `pandas`, в большинстве случаев достаточно вбить в поисковую строку "pandas/sklearn + (название вещи, про которую вы хотите узнать)".

GitHub страница библиотеки `featexp` с документацией (<https://github.com/abhayspawar/featexp>).

Отличный курс (<https://github.com/Yorko/mlcourse.ai>) по машинному обучению от OpenDataScience (сейчас он идет на английском, но русские итерации тоже бывают)

Задача 3.0.3. О чём эта новость? (60 баллов)

В этой задаче вам предстоит для некоторого набора новостных заголовков и текстов определить тематику новости.

Новость может быть одного из 7 типов:

1. Мир – “world”
2. Наука – “science”
3. Культура – “culture”
4. Экономика – “economy”
5. Политика – “politics”
6. Общество – “society”
7. Религия – “religion”

Особенностью данной задачи является то, что обучающий датасет и модель вам предстоит сделать самим.

Для сбора и структурирования датасета новостей жюри рекомендует использовать язык программирования Python и фреймворк Scrapy. Небольшой tutorial по этой библиотеке вы можете найти на официальном сайте фреймворка (<https://doc.scrapy.org/en/latest/intro/tutorial.html>).

Формат входных данных

Вам дан датасет с новостями (`test.csv`), в котором содержатся заголовки и тексты каждой статьи из подборки жюри (все эти новости реальные и взяты из архивов новостных сайтов разных годов).

Формат выходных данных

Вы должны сдать файл, в котором для каждой новостной статьи указана одна из тематик из списка выше. Правильный формат сдаваемого файла вы можете увидеть в файле (`sample_submission.csv` <https://drive.google.com/file/d/1EX5eUmNjgr5E3YpjCuZ6dmfohVIh1Izq/view?usp=sharing>).

Система оценки

Всего за эту задачу вы можете набрать до 60 баллов.

Определим `accuracy_score` как долю новостей, которые вы классифицировали правильно к общему количеству новостей в `test.csv`. Чем эта метрика выше, тем больше баллов вы получите, а именно:

1. Если $0.3 \leq \text{accuracy_score} \leq 0.5$, то вы получаете 10 баллов
2. Если $0.5 \leq \text{accuracy_score} \leq 0.6$, то вы получаете 20 баллов.
3. Если $0.6 \leq \text{accuracy_score} \leq 0.7$, то вы получаете 30 баллов.
4. Если $0.7 \leq \text{accuracy_score} \leq 0.8$, то вы получаете 40 баллов.
5. Если $0.8 \leq \text{accuracy_score} \leq 0.9$, то вы получаете 50 баллов.
6. Если $0.9 \leq \text{accuracy_score} \leq 1.0$, то вы получаете 60 баллов.

Решение

```

1  # Флаг, определяет, использовать ли видеокарту при подсчёте
2  USE_CUDA = True
3
4  import pandas as pd
5  import numpy as np
6  import time
7
8  ### Загрузим данные небольшой порцией
9  #data_text = pd.read_json('./news2.jsonlines', lines=True, chunksize=1000,
10 #encoding='utf8')
11 # data_text = pd.read_csv('./news2.csv', chunksize=1000, encoding='utf8')
12 # data_text = next(iter(data_text))
13
14 ### Для загрузки всего датасета в оперативную память (её может не хватить)
15 data_text = pd.read_csv('./news2.csv', encoding='utf8')
16
17 from sklearn.preprocessing import LabelEncoder
18
19 label_enc = LabelEncoder()
20 data_text['class'] = label_enc.fit_transform(data_text['class'])

```

Загруженные данные разделим на train и test, сохраним их, т.к. DataLoader'ы считывают данные из файлов

```

1  from sklearn.model_selection import train_test_split
2  data_train, data_val = train_test_split(data_text, test_size=0.5)
3
4  data_train.to_csv('train.csv', encoding='utf-8', index=False)
5  data_val.to_csv('val.csv', encoding='utf-8', index=False)
6
7  import torch
8  import torch.nn as nn
9  import torch.optim as optim
10 from torch.optim import lr_scheduler

```

rumorphy2 - библиотека, схожая с NLTK, но предназначенная для русского языка

morphy_tokenizer - лемматизирует слова в предложении и возвращает их список

Лемматизация — процесс приведения словоформы к лемме — её нормальной (словарной) форме.

```

1  import rumorphy2
2  morph = rumorphy2.MorphAnalyzer()
3  def morphy_tokenizer(s):
4      import re
5      lemmatized_words = []
6
7      for word in re.split("[^а-яА-ЯёЁА-За-я]", s):
8          word_repr = morph.parse(word.replace('.', ''))
9          if len(word_repr) != 0:
10             normal_form = word_repr[0].lexeme[0].word
11             lemmatized_words.append(normal_form)
12  return lemmatized_words

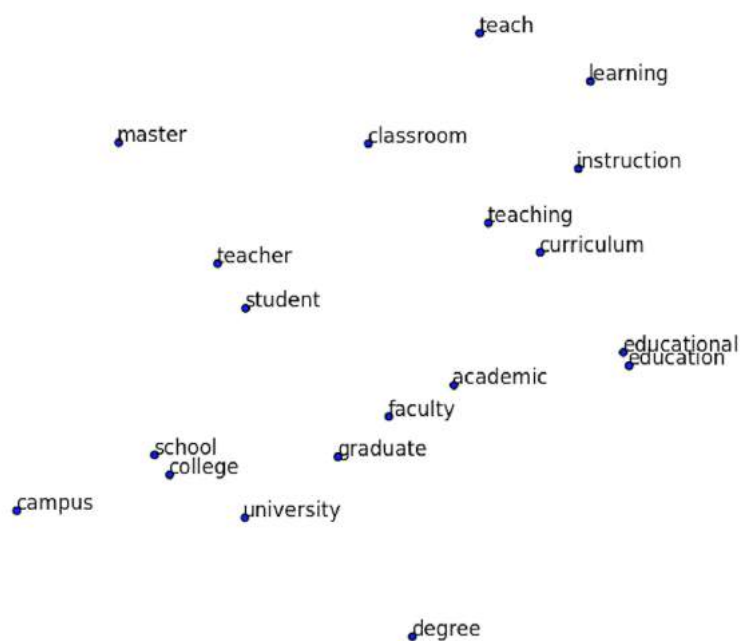
```

Создадим класс рекуррентного слоя

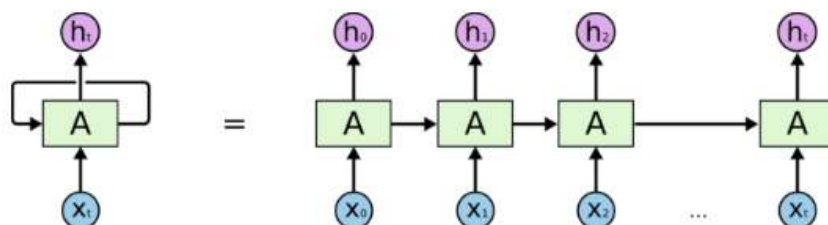
Embedding - замена каждого слова на вектор из некоторого количества вещественных чисел.

Indices	Latent Factors					
1	.32	.02	.48	.21	.56	.15
2	.65	.23	.41	.57	.03	.92
3	.45	.87	.89	.45	.12	.01
4	.65	.21	.25	.45	.78	.82

Так может выглядеть embedding матрица, в которой каждому индексу сопоставляется несколько чисел. Если смотреть embedding вектора как на вектора в некотором пространстве, то окажется, что близкие по смыслу слова окажутся близко и в этом пространстве.

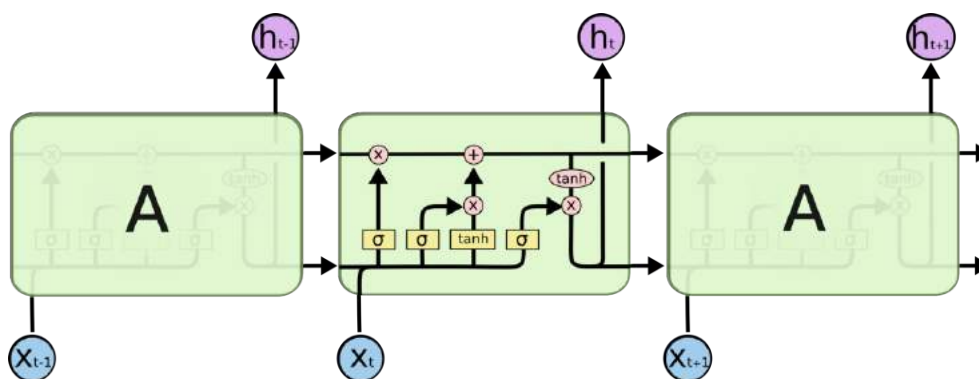


Сверху двумерные вектора слов после embedding'a.



An unrolled recurrent neural network.

Обычная рекуррентная нейросеть.



LSTM ячейка. Благодаря начилию долгосрочной памяти, она лучше обрабатывает длинные последовательности.

```

1 from torch.nn.utils.rnn import pack_padded_sequence, pad_packed_sequence
2
3 class RNN(nn.Module):
4
5     def __init__(self, hidden_size, encoder, num_layers=1):
6         '''
7         Args:
8             hidden_size: размер вектора внутреннего состояния LSTM
9             encoder: ссылка на слой эмбединга
10            num_layers: количество LSTM слоев
11        '''
12        super(RNN, self).__init__()
13        # embedding
14        self.encoder = encoder
15        # dropout после embedding
16        self.drop_en = nn.Dropout(p=0.6)
17        # рекуррентный слой
18        self.rnn = nn.LSTM(input_size=encoder.weight.size(1),
19                            hidden_size=hidden_size, num_layers=num_layers, dropout=0.5,
20                            batch_first=True, bidirectional=True)
21        # batchnorm после рекуррентных слоев
22        self.bn2 = nn.BatchNorm1d(hidden_size*2)
23
24    def forward(self, x, seq_lengths):
25        '''
26        Args:
27            x: (batch, time_step, input_size)
28            seq_lengths: длины последовательностей в x (чтобы их можно было
29                положить в тензор, они дополнены нулями справа)
30        Returns:
31            Выход полученный на последнем полносвязном слое без применения Softmax.
32        '''
33        # Применяем embedding и dropout
34        x_embed = self.encoder(x)
35        x_embed = self.drop_en(x_embed)
36        # Запаковываем последовательности, чтобы рекуррентный слой мог работать с
37        # ними, несмотря на то, что они имеют разные длины
38        packed_input = pack_padded_sequence(x_embed, seq_lengths.cpu().numpy(),
39                                           batch_first=True)
40
41        # Получаем выход с рекуррентного слоя
42        packed_output, ht = self.rnn(packed_input, None)
43        # Распаковываем полученный выход
44        out_rnn, _ = pad_packed_sequence(packed_output, batch_first=True)

```

```

45
46     row_indices = torch.arange(0, x.size(0)).long()
47     col_indices = seq_lengths - 1
48
49     if USE_CUDA:
50         row_indices = row_indices.cuda()
51         col_indices = col_indices.cuda()
52
53     # Из всех выходов LSTM получаем только выход последнего ненулевого символа
54     last_tensor=out_rnn[row_indices, col_indices, :]
55     # Применяем BatchNorm
56     out = self.bn2(last_tensor)
57
58     return out

```

Класс рекуррентной нейросети

Зададим архитектуру нейронной сети: у неё будет два входа для заголовка и текста. Для кодирования слов здесь используем word embedding

```

1  class Model(nn.Module):
2      def __init__(self, embedding_matrix, lstm_hdn, num_classes):
3          '''
4          Args:
5              embedding_matrix: матрица предобученных embedding'ов
6              lstm_hdn: размер внутреннего состояния для LSTM
7              num_classes: размерность вектора-результата (т.е. количество классов,
8                  на которые надо классифицировать)
9          '''
10         super().__init__()
11         # Создаем embedding слой и выключаем для него градиенты
12         self.encoder = nn.Embedding(embedding_matrix.size(0), embedding_matrix.size(1),
13             _weight=embedding_matrix)
14         self.encoder.weight.requires_grad = False
15
16         # предобученный embedding задает размер
17         self.text_rnn = RNN(lstm_hdn, self.encoder, num_layers=1)
18         self.title_rnn = RNN(lstm_hdn, self.encoder, num_layers=1)
19
20         # Выходной слой
21         self.l_out = nn.Linear(lstm_hdn*4, len(label_enc.classes_))
22
23     def forward(self, X_title, X_text, title_lengths, text_lengths,
24         title_perm, text_perm):
25         # Так как для LSTM слоя важно, чтобы длины последовательностей в батче
26         # шли в невозрастающем порядке, наш dataloader возвращает перестановку
27         # (permutation), которая упорядочивает последовательности.
28         # В рекуррентную часть мы передаем последовательности отсортировано,
29         # а затем возвращаем исходный порядок
30         title_outp = self.title_rnn(X_title[title_perm], title_lengths[title_perm])
31         title_outp[title_perm] = title_outp
32
33         # Делаем то же самое для текста
34         text_outp = self.text_rnn(X_text[text_perm], text_lengths[text_perm])
35         text_outp[text_perm] = text_outp
36
37         # Объединяем выходы с двух рекуррентных модулей и отправляем в полносвязный слой
38         return self.l_out(torch.cat([title_outp, text_outp], dim=1))

```

Создадим класс-загрузчик данных

Класс-генератор загрузит данные из файла и преобразует их морфизатором и токенайзером. А также приведёт данные к виду, пригодному для использования в рекуррентной нейросети

```

1 class TextClassDataLoader(object):
2
3 def __init__(self, path_file, word_to_index, batch_size=32, shuffle=True):
4     """
5     Args:
6         path_file:
7         word_to_index:
8         batch_size:
9     """
10
11 self.batch_size = batch_size
12 self.word_to_index = word_to_index
13 self.shuffle = shuffle
14
15 if word_to_index is None:
16     self.word_to_index = {'null': 0}
17
18 # read file
19 df = pd.read_csv(path_file).fillna('')
20 df['text'] = df['text'].apply(morphy_tokenizer)
21 df['text'] = df['text'].apply(self.generate_indexifyer())
22
23 print("Columns: {}".format(list(df.columns.values)))
24
25 df['title'] = df['title'].apply(morphy_tokenizer)
26 df['title'] = df['title'].apply(self.generate_indexifyer())
27 self.samples = df[['title', 'text', 'class']].values.tolist()
28
29 # for batch
30 self.n_samples = len(self.samples)
31 self.n_batches = int(self.n_samples / self.batch_size)
32 self.max_length = self._get_max_length()
33 self._shuffle_indices()
34
35 self.report()
36
37 def _shuffle_indices(self):
38     if not self.shuffle:
39         self.indices = np.arange(0, self.n_samples, 1)
40     else:
41         self.indices = np.random.permutation(self.n_samples)
42     self.index = 0
43     self.batch_index = 0
44
45 def _get_max_length(self):
46     length = 0
47     for sample in self.samples:
48         length = max(length, len(sample[1]))
49     return length
50
51 def generate_indexifyer(self):
52     def indexify(lst_text):
53         indices = []
54         for word in lst_text:
55             if word in self.word_to_index:
56                 indices.append(self.word_to_index[word])

```

```

57         else:
58             indices.append(0)
59         return indices
60
61     return indexify
62
63     @staticmethod
64     def _padding(batch_x):
65         batch_s = sorted(batch_x, key=lambda x: len(x))
66         size = len(batch_s[-1])
67         for i, x in enumerate(batch_x):
68             missing = size - len(x)
69             batch_x[i] = batch_x[i] + [0 for _ in range(missing)]
70         return batch_x
71
72     def _create_batch(self):
73         batch = []
74         n = 0
75         while n < self.batch_size:
76             _index = self.indices[self.index]
77             batch.append(self.samples[_index])
78             self.index += 1
79             n += 1
80         self.batch_index += 1
81
82         title, text, label = tuple(zip(*batch))
83
84         # получаем длины всех текстов в батче
85         text_seq_lengths = torch.LongTensor(list(map(len, text)))
86
87         # создаем тензор, в котором в каждой строчке находится последовательность,
88         # дополненная нулями справа.
89         # NOTE: длина строки тензора должна быть такой же, как и самая длинная
90         # последовательность мы могли бы взять большую длину, но в этом нет смысла,
91         # так как все последовательности уже точно поместятся
92         text_seq_tensor = torch.zeros((len(text), text_seq_lengths.max())).long()
93         for idx, (seq, seqlen) in enumerate(zip(text, text_seq_lengths)):
94             text_seq_tensor[idx, :seqlen] = torch.LongTensor(seq)
95
96         # Прodelывем то же самое для заголовков
97         title_seq_lengths = torch.LongTensor(list(map(len, title)))
98
99         # создаем тензор, в котором в каждой строчке находится последовательность,
100        # дополненная нулями справа.
101        # NOTE: длина строки тензора должна быть такой же, как и самая длинная
102        # последовательность мы могли бы взять большую длину, но в этом нет смысла,
103        # так как все последовательности уже точно поместятся
104        title_seq_tensor = torch.zeros((len(title), title_seq_lengths.max())).long()
105        for idx, (seq, seqlen) in enumerate(zip(title, title_seq_lengths)):
106            title_seq_tensor[idx, :seqlen] = torch.LongTensor(seq)
107
108        # Запомним перестановку, которая делает последовательности упорядоченными по длине,
109        # чтобы передать ее в саму модель. Это важно так как в стандартный модуль LSTM
110        # можно передавать только последовательности, идущие по уменьшению длины.
111        _, text_perm_idx = text_seq_lengths.sort(0, descending=True)
112        _, title_perm_idx = title_seq_lengths.sort(0, descending=True)
113
114        # Переведем лэйбл в тензор, чтобы не пришлось делать этого потом.
115        # (Для работы функции потерь лэйблы должны лежать в pytorch тензоре)
116        label = torch.LongTensor(label)

```

```

117
118     return text_seq_tensor, title_seq_tensor, label, text_seq_lengths, \
119            title_seq_lengths, text_perm_idx, title_perm_idx
120
121 def __len__(self):
122     return self.n_batches
123
124 def __iter__(self):
125     self._shuffle_indices()
126     for i in range(self.n_batches):
127         if self.batch_index == self.n_batches:
128             raise StopIteration()
129         yield self._create_batch()
130
131 def show_samples(self, n=10):
132     for sample in self.samples[:n]:
133         print(sample)
134
135 def report(self):
136     print('# samples: {}'.format(len(self.samples)))
137     print('max len: {}'.format(self.max_length))
138     print('# vocab: {}'.format(len(self.word_to_index)))
139     print('# batches: {} (batch_size = {})'.format(self.n_batches, self.batch_size))

```

Word embedding

Для эмбединга будем использовать часть предобученной на Википедии эмбединг-матрицы из fasttext. Она сопоставит каждому слову специальный вектор из 300 чисел.

<https://fasttext.cc/docs/en/crawl-vectors.html>

(файл прилагается с тетрадкой)

```

1 import pickle
2 word_to_index = pickle.load(open('word_index.pkl', 'rb'))
3 emb_tensor = torch.FloatTensor(np.load('embedding_matrix.npy', encoding='bytes'))
4
5 batch_size = 32
6
7 dataloaders = {'train': TextClassDataLoader('./train.csv', word_to_index, batch_size),
8               'val': TextClassDataLoader('./val.csv', word_to_index, batch_size)}
9 dataset_sizes = {'train': dataloaders['train'].n_samples,
10                 'val': dataloaders['val'].n_samples}
11
12 from tqdm import tqdm_notebook
13 def train_model(model, criterion, optimizer, scheduler, num_epochs=25):
14     since = time.time()
15
16     for epoch in range(num_epochs):
17         print('Epoch {}/{}'.format(epoch, num_epochs - 1))
18         print('-' * 10)
19
20         # каждая эпоха имеет обучающую и тестовую стадии
21         for phase in ['train', 'val']:
22             if phase == 'train':
23                 scheduler.step()
24                 model.train(True) # установить модель в режим обучения
25             else:
26                 model.train(False) # установить модель в режим предсказания

```



```

16
17     running_loss = 0.0
18     running_corrects = 0
19
20     # итерируемся по батчам
21     for data in tqdm_notebook(dataloaders[phase]):
22         new_data = list(data)
23         if USE_CUDA:
24             for i in range(len(data)):
25                 new_data[i] = data[i].cuda()
26
27         # получаем картинки и метки
28         X_title, X_text, label, title_lengths, text_lengths,
29             title_perm, text_perm = new_data
30
31         # инициализируем градиенты параметров
32         optimizer.zero_grad()
33
34         # forward pass
35         outputs = model(X_title, X_text, title_lengths, text_lengths,
36                         title_perm, text_perm)
37         _, preds = torch.max(outputs.data, 1)
38         loss = criterion(outputs, label)
39
40         # backward pass + оптимизируем только если это стадия обучения
41         if phase == 'train':
42             loss.backward()
43             optimizer.step()
44
45         # статистика
46         running_loss += loss.item()
47         running_corrects += int(torch.sum(preds == label.data))
48
49         # Выводим статистику о качестве модели во время обучения
50         epoch_loss = running_loss / dataset_sizes[phase]
51         epoch_acc = running_corrects / dataset_sizes[phase]
52
53         print('{ } Loss: {:.4f} Acc: {:.4f}'.format(
54             phase, epoch_loss, epoch_acc))
55
56     print()
57
58     time_elapsed = time.time() - since
59     print('Training complete in {:.0f}m {:.0f}s'.format(
60         time_elapsed // 60, time_elapsed % 60))
61
62     return model, losses

```

Обучение

Передадим в нашу модель embedding матрицу

```

1 model = Model(emb_tensor, 32, len(label_enc.classes_))
2 if USE_CUDA:
3     model = model.cuda()
4
5 # В качестве cost function используем кросс-энтропию
6 loss_fn = nn.CrossEntropyLoss()
7
8 # В качестве оптимизатора - стохастический градиентный спуск

```

```

9 optimizer_ft = optim.Adam(model.parameters(), lr=0.001)
10
11 # Умножает learning_rate на 0.1 каждые 7 эпох
12 # (это одна из эвристик, не было на лекциях)
13 exp_lr_scheduler = lr_scheduler.StepLR(optimizer_ft, step_size=7, gamma=0.1)

```

Обучение происходит здесь

Процесс необходимо остановить, как только val Loss перестанет убывать

```

1 train_model(model, loss_fn, optimizer_ft, exp_lr_scheduler, num_epochs=25)

```

Сохраним модель

```

1 torch.save(model.state_dict(), './model.pth')
2
3 model.load_state_dict(torch.load('./model.pth'))

```

Предсказание на тестовой выборке

```

1 data_test = pd.read_csv('FINAL_TEST.csv')
2 data_test['class'] = 0
3 data_test.to_csv('./test_with_dummy_class.csv', index=False)
4 data_test
5
6 # Создадим массив из нулей для будущих предсказаний
7 predictions = np.zeros((762,))
8 index_to_write = 0
9
10 model.train(False)
11 for data in tqdm_notebook(TextClassDataLoader('./test_with_dummy_class.csv',
12                                             word_to_index, 32, shuffle=False)):
13     new_data = list(data)
14     if USE_CUDA:
15         for i in range(len(data)):
16             new_data[i] = data[i].cuda()
17
18     # получаем картинки и метки
19     X_title, X_text, label, title_lengths, text_lengths, title_perm,
20     text_perm = new_data
21
22     # forward pass
23     outputs = model(X_title, X_text, title_lengths, text_lengths, title_perm, text_perm)
24     _, preds = torch.max(outputs.data, 1)
25
26     predictions[index_to_write: min(index_to_write+32, 762)] = preds.cpu().numpy()
27     index_to_write += 32
28
29 predictions = pd.DataFrame({'category': predictions})
30 predictions.category = predictions.category.map({i: class_name for i, class_name
31         in enumerate(label_enc.classes_)})
32 predictions.to_csv('res.csv', encoding='utf-8', index=False)

```

Основы для кода RNN модуля и dataloader'a взяты из репозитория Pytorch-RNN-text-classification. (<https://github.com/keishinkickback/Pytorch-RNN-text-classification>)

Статья на русском про LSTM. (<http://datareview.info/article/issleduem-lstm-seti-chast-1/>)

Чтобы быстро учить нейросеть вы можете использовать Google Collab или, если у вас есть видеокарта с поддержкой CUDA, установить CUDA Toolkit и PyTorch с поддержкой CUDA на свой компьютер. Небольшое руководство (<https://habr.com/post/348058/>) по использованию GoogleCollab (в нем используется TensorFlow, но разница с PyTorch незначительна). Руководство по настройке CUDA (<https://docs.nvidia.com/cuda/cuda-installation-guide-microsoft-windows/index.html>) на официальном сайте достаточно полное, но у многих возникают ошибки, решение которых придется искать самому.