

Большие данные и машинное обучение

Инженерный тур

Легенда задачи

- Есть лента «новостей» в некоторой социальной сети, которая наполнена различным контентом.
- Контент состоит из множества записей, которые показываются пользователю.
- Пользователь просматривает каждую запись и как-то взаимодействует ней.
- Для построения системы рекомендации требуется предсказывать тип взаимодействия.

Описание оборудования и программного обеспечения

Участникам предоставляется доступ к виртуальному серверу для проведения вычислений:

- 12 ГБ ОЗУ;
- 40 Гб жесткий диск;
- 4 ядра процессора;
- ОС Ubuntu 16.04+.

На сервере установлен Python 3.9 с Jupyter Notebook, SciPy, NumPy, Scikit Learn и Pandas. Участники могут доустановить необходимые им библиотеки и языки.

Доступ к серверу осуществляется с личного или выданного ноутбука через веб-интерфейс Jupyter Notebook.

Решения загружаются на платформе Codeforces.

Описание этапов работы участников

1. Знакомство с задачей.
2. Разработка и реализация подхода к работе с данными, которые хранятся в нескольких таблицах.
3. Разработка программы, которая построит модель предсказания.
4. Поиск оптимальных гиперпараметров выбранной модели предсказания.
5. Построение предсказания на тестовом множестве и загрузка его в тестирующую систему.

Задача VI.2.3.1. Предсказание оценки (100 баллов)

Имя входного файла: стандартный ввод.

Имя выходного файла: стандартный вывод.

Ограничение по времени выполнения программы: 15 с.

Ограничение по памяти: 256 Мбайт.

Условие

Требуется предсказать, понравится ли пользователям запись в социальной сети?

Формат входных данных

Набор данных, подготовленный командой VK, доступен по адресу <https://disk.yandex.ru/d/Gffq-VhB5jekKQ>.

Он состоит из нескольких CSV файлов:

- `topics.csv` — содержит информацию о записях социальной сети «Одноклассники». Каждая запись характеризуется признаками текста и изображения.
- `users.csv` — содержит информацию о пользователях социальной сети «Одноклассники». Каждый пользователь характеризуется датой рождения, полом и идентификатором города.
- `train.csv` — содержит информацию о взаимодействии пользователей с записями социальной сети «Одноклассники». Тип взаимодействия равен L, если запись понравилась, и D — если запись не понравилась.
- `test.csv` — содержит идентификаторы записей и пользователей, для которых требуется предсказать тип взаимодействия.

Формат выходных данных

Необходимо загрузить в тестирующую систему текстовый файл с ответами. Он не должен содержать заголовки и никакой другой дополнительной информации. Каждая строка должна содержать одну заглавную латинскую букву D или L. Ответы должны быть даны в порядке, в котором были соответствующие запросы из файла `test.csv`.

Вы можете совершить **не более 50-ти попыток** решения данной задачи.

Критерии оценивания

Для оценки будет использоваться F-мера класса L, умноженная на 10000.

Результат, который вы будете получать во время соревнования, будет вычислен только на 20% тестовой выборки. В конце соревнования выбранные вами решения **будут перетестированы** на оставшихся 80%.

Решение

Задачу можно решить несколькими способами.

В качестве базового решения используется подход, в котором для каждой записи вычисляется число отметок типа L (нравится) и D (не нравится). Если отметок L было больше, то для всех запросов из тестового множества будет даваться ответ L независимо от пользователя, иначе будет даваться ответ D.

Можно решить эту задачу как классическую задачу обучения с учителем. Для этого требуется объединить таблицы с описанием пользователей и записей. Но в результате будут получаться объекты из большого числа признаков. Для построения модели можно воспользоваться библиотекой XGBoost.

Также можно решить эту задачу как классическую задачу коллаборативной фильтрации. Для этого достаточно использовать только таблицу с информацией о взаимодействиях пользователей и записей. Можно представить, что эта таблица кодирует разреженную матрицу размера N на M , где N — число пользователей, а M — число записей. Взаимодействие L можно кодировать значением $+1$, а D — значением -1 . В таком виде можно попытаться факторизовать матрицу на две матрицы с небольшим рангом. Например, для этого можно использовать метод NMF из библиотеки `surprise`.

Помимо этого можно объединить предыдущие два решения, если сконкатенировать полученные при разложении матриц признаки с данным признаками пользователей и записей.

Пример программы-решения

Ниже представлено решение на языке Python 3.

```
1 import pandas as pd
2 import numpy as np
3 import pickle
4 from sklearn.model_selection import KFold
5 from sklearn.metrics import f1_score
6 import seaborn as sns
7 import warnings
8 warnings.filterwarnings('ignore')
9 OPTIMIZE_ROUNDS = False
10
11 from catboost import cv, Pool, CatBoostClassifier
12
13 train = pd.read_csv('/kaggle/input/nonononono/trainavrdd.csv')
14 test = pd.read_csv('/kaggle/input/nonononono/testavrdd.csv')
15
16 train.birthDate.fillna(train.birthDate.median(), inplace = True)
17 train.drop('Unnamed: 0', axis = 1, inplace = True)
18 test.drop('Unnamed: 0', axis = 1, inplace = True)
19
20 train['locationId']=train['locationId'].astype("category")
21 train['gender']=train['gender'].astype("category")
22 train['topicId']=train['topicId'].astype("category")
23 train['userId']=train['userId'].astype("category")
24 test['locationId']=test['locationId'].astype("category")
25 test['gender']=test['gender'].astype("category")
26 test['topicId']=test['topicId'].astype("category")
27 test['userId']=test['userId'].astype("category")
28
29 train['birthDate'] = (train.birthDate - train.birthDate.min()) /
    ↪ (train.birthDate.max() - train.birthDate.min())
30 train['countus'] = (train.countus - train.countus.min()) / (train.countus.max() -
    ↪ train.countus.min())
31 train['counttop'] = (train.counttop - train.counttop.min()) /
    ↪ (train.counttop.max() - train.counttop.min())
32
33 test.rename({'countus_y':'countus'}, axis=1)
34
```

```

35 test['birthDate'] = (test.birthDate - test.birthDate.min()) /
    ↪ (test.birthDate.max() - test.birthDate.min())
36 test['countus_y'] = (test.countus_y - test.countus_y.min()) /
    ↪ (test.countus_y.max() - test.countus_y.min())
37 test['counttop'] = (test.counttop - test.counttop.min()) / (test.counttop.max() -
    ↪ test.counttop.min())
38
39 train.drop(['Unnamed: 0.1', 'locationId', 'userId', 'topicId'], axis=1,
    ↪ inplace=True)
40
41 y = train['feedback']
42 X = train.drop('feedback', axis = 1)
43
44 y = y.replace([-1], 0)
45
46 K = 5
47 kf = KFold(n_splits = K, random_state = 1, shuffle = True)
48
49 model = CatBoostClassifier(
50     cat_features = ['gender'],
51     depth = 16,
52     l2_leaf_reg = 34,
53     random_seed=42,
54     learning_rate=0.00009,
55     iterations = 50,
56     loss_function='Logloss'
57 )
58
59 for i, (train_index, test_index) in enumerate(kf.split(train)):
60     y_train, y_valid = y.iloc[train_index], y.iloc[test_index]
61     X_train, X_valid = X.iloc[train_index,:], X.iloc[test_index,:]
62     print( "\nFold ", i)
63     if OPTIMIZE_ROUNDS:
64         fit_model = model.fit( X_train, y_train,
65                               eval_set=[X_valid, y_valid],
66                               use_best_model=True
67         )
68
69     else:
70         fit_model = model.fit( X_train, y_train )
71
72 test.drop(['Unnamed: 0.1', 'locationId', 'userId', 'topicId'],
    ↪ axis=1,inplace=True)
73 model.get_feature_importance(prettified=True)
74 test.drop(['Unnamed: 0.1', 'locationId', 'userId', 'topicId'],
    ↪ axis=1,inplace=True)
75
76 f=(model.predict(test))
77
78 d1=[]
79 for i in f:
80     if i == 1:
81         d1.append('L')
82     else:
83         d1.append('D')
84 len(d1)
85
86 g1=True
87 for i in d1:
88     if g1==True:

```

```
89     my_file = open("ff12.txt", "w+")
90     my_file.write(i)
91     my_file.write('\n')
92     my_file.close()
93     g1=False
94     else:
95         my_file = open("ff12.txt", "a+")
96         my_file.write(i)
97         my_file.write('\n')
98         my_file.close()
```

Материалы для подготовки

- <https://github.com/deeppavlov/dlschl>.
- <https://stepik.org/course/1296/promo>.
- <https://stepik.org/course/217/promo>.
- <https://stepik.org/course/76/promo>.
- https://drive.google.com/file/d/1i1_Wf8jiWd3o0nRZGaX2gdoHCTIkvK0J/view.