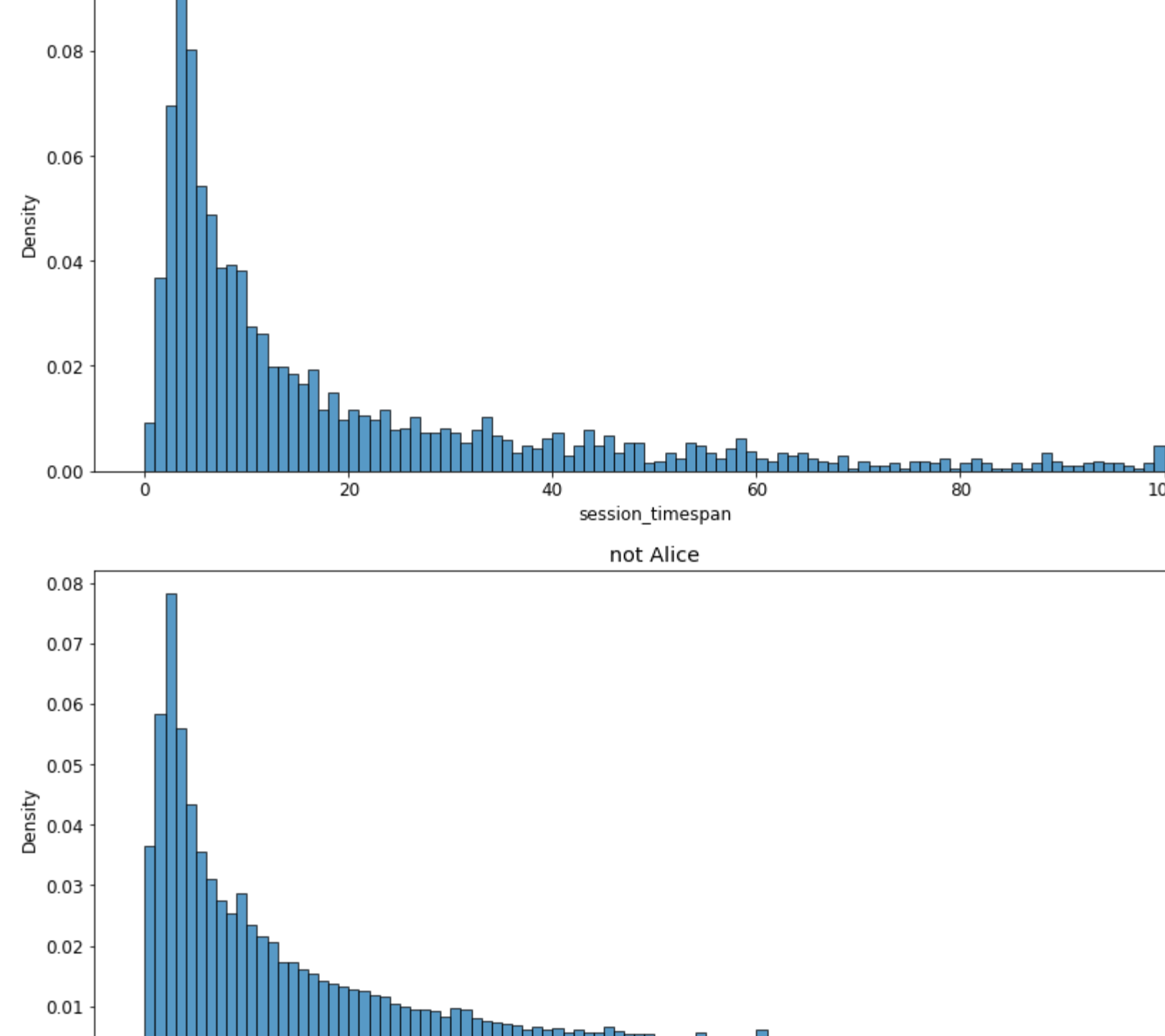


Учебный проект  
предоставленные  
разделы:



Alice



Вблизи нулевых значений продолжительности сессии, гистограммы заметно отличаются, вероятно это то, что нам надо. Можно построить хитрую тепловую карту, в какой ячейке которой будет содержаться разница между процентным содержанием сессии данной продолжительности у Alice и у всех, кроме Alice.

Например, если с продолжительности 0 секунд в наборе сессий у Alice встречается в 0.83% всех сессий, а у всех, кроме Alice в 2.61%, значит, разница для сессии продолжительности 0 секунд равна  $1.78\%$ . Это может, например, означать, что Alice меньше других склонна открывать много вкладок за раз.

При необходимости в определённую ниже функцию `per_difference_heatmap` можно передавать количество секунд для группировки.

```
In (36): def per_difference_heatmap(seconds_to_group=1):
    """Создаёт heatmap для разницы продолжительности сессии (Alice - not Alice)
    Продолжительные значения означают, что сессии Alice чаще других имеют данную продолжительность,
    продолжительности. При необходимости можно изменять временное окно (seconds_to_group)"""
    start_second = 0
    stop_second = 100 # seconds to group

    def session_timespan_info(df, bins):
        times = df['session_timespan'].values
        result = []
        for i in range(len(bins)-1):
            for j in range(1, len(times) - times[i] + 1):
                if times[i+j] >= bins[i] & times[i+j+1] < bins[i+1]:
                    result.append(j)
        return np.array(result)

    bins = [i for i in range(start_second, stop_second, seconds_to_group)]
    times_alice = session_timespan_info(train_df[train_df['target'] == 1], bins)
    times_not_alice = session_timespan_info(train_df[train_df['target'] == 0], bins)
    difference = times_alice - times_not_alice
    difference = difference / (difference.max() + difference.min())
    difference = difference * 100

    print(f'Минимальная разница: {round(difference.min(), 1)} %')
    print(f'Максимальная разница: {round(difference.max(), 1)} %')
    print(f'Возможные значения означают, что сессии Alice чаще других имеют данную продолжительность')
    labels = ['seconds_to_group']
    for sec, diff in enumerate(difference[16]).reshape(-1, 11):
        sns.heatmap(data=diff.reshape(-1, 11),
                    annot=True,
                    annot_knames=['seconds_to_group'],
                    yticklabels=False,
                    xticklabels=False,
                    cbar=False,
                    fmt='')

    per_difference_heatmap(seconds_to_group)
```

Получившиеся значения означают, что сессии Alice чаще других имеют данную продолжительность



На основе данной таблицы в дальнейшем будет построен признак, учитывающий особенности продолжительности сессии Alice.

## Признаки, связанные с набором сайтов в сессии

Следующая группа признаков, которую можно записать связана с набором сайтов в сессии. Подавляющее большинство признаков этой группы будет порождено в дальнейшем в виде матрицы частот с помощью `CountVectorizer` (или `TfidfVectorizer`), здесь же будут рассмотрены остальные признаки, которые могут быть извлечены из наборов сайтов в сессии.

1. Количество уникальных сайтов в сессии

Возможно, отличить Alice от остальных, количество различных сайтов, которые она использует в пределах одной сессии.

```
In (37): num_unique_sites = train_df['#unique_sites'].values

In (38): plt.subplots(1, 2, figsize = (15, 6))

ax = sns.countplot(num_unique_sites[y_train == 1])
ax.set(title='Alice',
        xlabel='#unique_sites')

plt.subplot(1, 2, 2)
ax = sns.countplot(num_unique_sites[y_train == 0])
ax.set(title='Not Alice',
        xlabel='#unique_sites')
```



Диаграмма для Alice имеет слабовазражаемые особенности, возможно, в дальнейшем стоит попробовать их использовать для построения признаков.

1. Сайты, встречающиеся чаще в сессии Alice, чем в сессиях других пользователей и наборот

В дальнейшем хотелось бы создать несколько признаков-индикаторов или счетчиков для определенных сайтов в сессии, наличие или количество которых могут быть дополнительно повышать качество классификации.

```
In (39): def attendance_list(df, target):
    """Возвращает топ сайтов (id сайтов, количество посещений)"""
    ind_cnt = np.unique(df[df['target'] == target])
    f1 = f1.groupby(ind_cnt).agg('count').reset_index()
    ind_cnt = np.sort(key=lambda x: x[1], reverse=True)
    return ind_cnt

In (40): def out_of_popular_sites(alice_first, not_alice_first, top_alice):
    """Возвращает множество сайтов, входящих в топ-Alice, но не входящих в топ-not Alice (top_alice=True),
    либо входящих в топ-not Alice, но не входящих в топ-Alice (top_alice=False)"""
    alice_top = attendance_list(train_df, target=1)[alice_first]
    not_alice_top = attendance_list(train_df, target=0)[not_alice_first]

    set_alice_top = set(alice_top)
    set_not_alice_top = set(not_alice_top)
    if top_alice is True:
        return set_alice_top - set_not_alice_top
    if top_alice is False:
        return set_not_alice_top - set_alice_top
```

for check = [(100, 3900), (60, 1000), (60, 2200), (20, 150), (10, 45), (10, 40), (5, 20)]

```
for alice_first, not_alice_first in for_check:
    print(f'Популярные сайты Alice, сайты, которые не входят в топ-not Alice (first): ')
    f1 = f1.groupby(ind_cnt).agg('count').reset_index()
    ind_cnt = np.sort(key=lambda x: x[1], reverse=True)
    return ind_cnt
```

Среди топ-100 сайтов Alice, сайты, которые не входят в топ-3900: (27307, 0), (2712, 0), (27189, 0), (25383, 0)  
Среди топ-60 сайтов Alice, сайты, которые не входят в топ-2200: (2080, 0), (12619, 0), (2307, 0)  
Среди топ-20 сайтов Alice, сайты, которые не входят в топ-150: (1300, 0)  
Среди топ-10 сайтов Alice, сайты, которые не входят в топ-45: (81, 0), (82, 0)  
Среди топ-5 сайтов Alice, сайты, которые не входят в топ-20: (81, 0), (82, 0), (89, 0)  
Среди топ-3 сайтов Alice, сайты, которые не входят в топ-5: (81, 0), (82, 0)

```
In (42): for check = [(100, 30), (100, 20), (100, 10), (100, 5)]

for alice_first, not_alice_first in for_check:
    print(f'Популярные сайты Alice, сайты, которые не входят в топ-not Alice (first): ')
    f1 = f1.groupby(ind_cnt).agg('count').reset_index()
    ind_cnt = np.sort(key=lambda x: x[1], reverse=True)
    return ind_cnt
```

Среди топ-100 сайтов, сайты, которые не входят в топ-3000: (27307, 0), (2712, 0), (27189, 0), (25383, 0)  
Среди топ-60 сайтов, сайты, которые не входят в топ-2200: (2080, 0), (12619, 0), (2307, 0)  
Среди топ-20 сайтов, сайты, которые не входят в топ-150: (1300, 0)  
Среди топ-10 сайтов, сайты, которые не входят в топ-45: (81, 0), (82, 0)  
Среди топ-5 сайтов, сайты, которые не входят в топ-20: (81, 0), (82, 0), (89, 0)  
Среди топ-3 сайтов, сайты, которые не входят в топ-5: (81, 0), (82, 0)

Заметим, что имеется ряд сайтов, которые Alice посещает чаще, чем остальные люди в наборе данных. Также имеются и такие сайты, которые Alice посещает реже, чем остальные люди. В дальнейшем будут сформированы два признака-индикатора входящего одного из интересующих сайтов в сессию.

Помогим на сайты, встречающиеся таким образом:

```
In (43): with open(os.path.join(PATH_TO_DICT, 'site_dict.pkl'), 'wb') as file:
    sites_dict = pickle.load(file)
    reversed_dict = {v: k for k, v in sites_dict.items()}
    reversed_dict[0] = '0'

In (44): sites_top_alice = [3000, 2080, 27307, 12619, 263, 25383, 81, 81, 879, 77]
sites_top_not_alice = [55, 56, 57, 778, 780, 782, 786, 812]

print("\n033[32m Alice (033[0m\n")
for site in sites_top_alice:
    print(f'{site}: {reversed_dict[site]}')

print("\n033[31m Not Alice (033[0m\n")
for site in sites_top_not_alice:
    print(f'{site}: {reversed_dict[site]}')
```

Alice

2080: media-1.melty.ru  
263: www.caf.ru  
25383: www.cdn.justice.gov.ru  
27307: fr.glee.wikia.com  
12619: yowach.ru  
77: 11.ytimg.com  
879: r1--sn-akobuk-jqbe.googlevideo.com  
81: r1--sn-akobuk-jqbe.googlevideo.com  
3000: vk.com

not Alice

778: www.ncbi.nlm.nih.gov  
780: blast.ncbi.nlm.nih.gov  
812: mail.google.com  
782: annotation.fr  
786: www.phylogeny.fr  
55: safebrowsering-cache.google.com  
56: safebrowsering-clients.google.com  
57: plus.google.com

Среди найденных сайтов имеются такие, которые в результате дополнительной проверки оказались неинформативными (Alice посещала их всего в нескольких сессиях) и в дальнейшем эти сайты (263, 25383) не использовались при построении признаков.

Видим, что Alice довольно часто использует французские сайты, возможно получится выделить подходящий признак.

1. Встречается ли домен "fr" в названии сайта

В предыдущем пункте было обнаружено, что в некоторых уникальных для Alice сайтов всегда встречается домен "fr". Следует сравнить частоту появления этого домена в сайтах Alice и в сайтах остальных пользователей.

```
In (45): def hist_for_counted_words(df, pattern, sites_dict, density=True):
    """Создаёт гистограммы для частот встречаемых сайтов в сессии,
    с подстройкой 'pattern'"""
    def search_and_sub(df, pattern):
        import copy
        new_df = copy.copy(df)
        sites = ['sites']
        for i in range(1, 11):
            new_df['sites'] = new_df['sites'].apply(lambda x: x if re.search(pattern, str(x)) else 0)
        return new_df

    print(f'Пикс: {pattern}')

    sites = ['sites']
    for i in range(1, 11):
        sites_df = df[sites].fillna(0).astype('int').applymap(lambda x: sites_dict[x])
        reversed_dict = {v: k for k, v in sites_dict.items()}
        reversed_dict[0] = '0'

    sites_df = search_and_sub(df, pattern)
    sites_df['count'] = sites_df[sites].values.sum(axis=1)

    plt.subplots(1, 2, figsize = (15, 6))

    plt.subplot(1, 2, 1)
    ax = sites_df[sites].hist(ranges=(1, 11), density=True)
    ax.set(title='Alice',
        xlabel='number of occurrences')

    plt.subplot(1, 2, 2)
    ax = sites_df[sites].hist(ranges=(1, 11), density=True)
    ax.set(title='Not Alice',
        xlabel='number of occurrences')
```

hist\_for\_counted\_words(train\_df, 'fr', reversed\_dict, density=False)



Существенных отличий двух распределений не наблюдается, что скорее всего означает, что признак будет неинформативным. К слову, кайли поддержку, по которой можно было бы построить признак, так и не представилось возможным.

В следующем разделе найденные особенности будут применены для добавления новых признаков к тренировочному и тестовому наборам данных.

## Обучение моделей (3)

В этом блоке используются файлы, полученные при помощи предыдущих блоков (тренировочный и тестовый наборы данных). Наборы векторизуются, к полученным признакам матрицы добавляются новые признаки, а затем на полученных данных обучается логистическая регрессия и применяется для классификации тестовых сессий.

### Содержание

1. Загрузка подготовленных данных
2. Применение `CountVectorizer` и `TfidfVectorizer`
3. Функции для добавления фичей
4. Применение преобразованных наборов данных для классификации
5. Применение новых фичей для классификации
6. Заключение

### Загрузка подготовленных данных

```
In (1): import os
import warnings
warnings.filterwarnings('ignore')

from tqdm import tqdm_notebook
import pickle

import numpy as np
import pandas as pd

import scipy.sparse as sp

from sklearn.model_selection import TimeSeriesSplit, cross_val_score, GridSearchCV
from sklearn.linear_model import LogisticClassifier, LogisticRegression, LogisticRegressionCV
from sklearn.metrics import roc_auc_score, accuracy_score
from sklearn.feature_extraction.text import CountVectorizer, TfidfVectorizer

PATH_TO_DATASET = os.path.join('intermediate_data', 'test_train')
PATH_TO_DICT = os.path.join('initial_data', 'site_dict')
```

```
def write_to_submission_file(predicted_labels, index_label='session_id'):
    predicted_df = pd.DataFrame(predicted_labels,
                                index=np.arange(1, predicted_labels.shape[0] + 1),
                                columns=[target])
    predicted_df.to_csv(PATH_TO_DICT, index_label=index_label)
```

```
In (2): with open(os.path.join(PATH_TO_DATASET, 'train_w10_w30_final.pkl'), 'rb') as f:
    train_df = pickle.load(f)

with open(os.path.join(PATH_TO_DATASET, 'test_w10_w30_final.pkl'), 'rb') as f:
    test_df = pickle.load(f)

train_df = train_df.sort_values(by=['year', 'month', 'day', 'time'])
```

train\_test\_df = pd.concat([train\_df, test\_df])  
train\_test\_df['site'] = train\_test\_df['site'].fillna(0).astype('int')  
y\_train = train\_test\_df['target'].astype('int').values

## Применение CountVectorizer и TfidfVectorizer

```
In (48): train_df['sites'] = [i for i in range(1, 11)]
sites = ['sites'].fillna(0).astype('int').to_csv(os.path.join('intermediate_data', 'train_sessions_text_cnt.txt'),
        sep=' ',
        index=None, header=None)
test_df['sites'] = [i for i in range(1, 11)]
test_df['sites'] = test_df['sites'].fillna(0).astype('int').to_csv(os.path.join('intermediate_data', 'test_sessions_text_cnt.txt'),
        sep=' ',
        index=None, header=None)

with open(os.path.join(PATH_TO_DICT, 'site_dict.pkl'), 'rb') as file:
    sites_dict = pickle.load(file)
    reversed_dict = {v: k for k, v in sites_dict.items()}
    reversed_dict[0] = '0'

train_df['sites'] = train_df['sites'].applymap(lambda x: reversed_dict[x].to_csv(os.path.join('intermediate_data', 'train_sessions_text_tfidf.txt'),
        sep=' ',
        index=None, header=None)
test_df['sites'] = test_df['sites'].applymap(lambda x: reversed_dict[x].to_csv(os.path.join('intermediate_data', 'test_sessions_text_tfidf.txt'),
        sep=' ',
        index=None, header=None)
```

```
In (49): #time
tfidf_vectorizer = CountVectorizer(ngram_range=(1, 3), max_features=50000, tokenizer=lambda s: s.split())
with open(os.path.join('intermediate_data', 'train_sessions_text_cnt.txt')) as inp_train_file:
    X_train = cnt_vectorizer.fit_transform(inp_train_file)
with open(os.path.join('intermediate_data', 'train_sessions_text_tfidf.txt')) as inp_train_file:
    X_train_tfidf = tfidf_vectorizer.fit_transform(inp_train_file)
with open(os.path.join('intermediate_data', 'test_sessions_text_cnt.txt')) as inp_test_file:
    X_test = cnt_vectorizer.transform(inp_test_file)
with open(os.path.join('intermediate_data', 'test_sessions_text_tfidf.txt')) as inp_test_file:
    X_test_tfidf = tfidf_vectorizer.transform(inp_test_file)
```

Wall time: 35.6 s

```
In (50): cnt_vectorizer.get_feature_names()[1:10]

Out(50): ['1086',
'11',
'110',
'111',
'1112',
'1114',
'1115',
'11121',
'1119',
'112']
```

tfidf\_vectorizer.get\_feature\_names()[1:10]

```
Out(51): ['0',
'0 0',
'0 0 0',
'0 0 0 0',
'0 academia-assets.com',
'0 docs.google.com',
'0 docs.google.com 0',
'0 docs.google.com 0',
'0 docs.google.com 0 0',
'0 docs.google.com 0 0 0']
```

## Функции для добавления фичей

```
In (8): def add_time_features(df, X_sparse):
    """Добавляет индикаторы (4 признака) утра, дня, вечера и ночи"""
    hour = df['start_hour']
    morning = (hour >= 7) & (hour <= 11).astype('int')
    day = (hour >= 12) & (hour <= 18).astype('int')
    evening = (hour >= 19) & (hour <= 23).astype('int')
    night = (hour >= 0) & (hour <= 6).astype('int')
    X = sp.hstack(X_sparse, morning.values.reshape(-1, 1),
        day.values.reshape(-1, 1),
        evening.values.reshape(-1, 1),
        night.values.reshape(-1, 1))
    return X
```

```
In (9): def add_session_timespan(df, X_sparse):
    """Добавляет продолжительность сессии в секундах (1 признак)"""
    session_timespan = df['session_timespan']
    X = sp.hstack(X_sparse, session_timespan.values.reshape(-1, 1))
    return X
```

```
In (10): def feature_short_session(df, X_sparse, time_min=0, time_max=3):
    """Добавляет индикаторы короткой сессии (от 0 до 2 секунд)
    (1 признак)"""
    def short_session(x):
        if (x >= time_min) and (x < time_max):
            return 1
        return 0
    times = df['session_timespan']
    X = sp.hstack(X_sparse, times.apply(short_session).values.reshape(-1, 1))
    return X
```

```
In (11): def feature_middle_session(df, X_sparse, time_min=3, time_max=9):
    """Добавляет индикаторы средней сессии (от 3 до 6 секунд)
    (1 признак)"""
    def middle_session(x):
        if (x >= time_min) and (x < time_max):
            return 1
        return 0
    times = df['session_timespan']
    X = sp.hstack(X_sparse, times.apply(middle_session).values.reshape(-1, 1))
    return X
```

```
In (12): def add_top_sites(df, X_sparse, list_of_sites=[3000, 2080, 27307]):
    """Добавляет индикаторы (1 признак) нескольких сайтов, которые Alice
    использует чаще всего, чем остальные люди"""
    sites = df['sites']
    result = [0] * len(sites)
    for i, site in enumerate(sites):
        if site in list_of_sites:
            result[i] = 1
    X = sp.hstack(X_sparse, np.array(result).reshape(-1, 1))
    return X
```

```
In (13): """Оказалось, что признак плохой и приводит к переобучению"""
def add_start_hour(df, X_sparse):
    """Добавляет час начала сессии (1 признак)"""
    start_hour = df['start_hour']
    X = sp.hstack(X_sparse, start_hour.values.reshape(-1, 1))
    return X
```

```
In (14): def add_day_of_week(df, X_sparse):
    """Добавляет день начала сессии (1 признак)"""
    day_of_week = df['day_of_week']
    X = sp.hstack(X_sparse, day_of_week.values.reshape(-1, 1))
    return X
```

Далее закомментированы определения функций, добавляющих плохие признаки.

```
In (15): """
def add_num_of_unique(df, X_sparse):
    """Добавляет количество уникальных сайтов в сессии (1 признак)"""
    num_of_unique = df['#unique_sites']
    X = sp.hstack(X_sparse, num_of_unique.values.reshape(-1, 1))
    return X
```

```
Out(15): """def add_num_of_unique(df, X_sparse):\n    # Добавляет количество уникальных сайтов в сессии (1 признак)\n    num_of_unique = df['#unique_sites']\n    X = sp.hstack(X_sparse, num_of_unique.values.reshape(-1, 1))\n    return X"""
```

```
In (16): """
from sklearn.preprocessing import StandardScaler
def add_session_timespan_scaled(df, X_sparse):
    X = sp.hstack(X_sparse, session_timespan.values.reshape(-1, 1))
    scaler = StandardScaler()
    scaler.fit(X_sparse)
    X_train_scaled = scaler.transform(X_sparse)
    X_test_scaled = scaler.transform(X_test)
    return X_train_scaled, X_test_scaled
```

```
Out(16): """from sklearn.preprocessing import StandardScaler\nX_train_scaled, X_test_scaled = StandardScaler().fit_transform(X_train)\nX_test_scaled, X_test_scaled = StandardScaler().fit_transform(X_test)\nreturn X_train_scaled, X_test_scaled\n"""
```

## Применение преобразованных наборов данных для классификации

В данном разделе были обучены `LogisticRegressionCV` на двух тренировочных наборах данных:

- К тренировочным и тестовым наборам, преобразованным при помощи `CountVectorizer` и `TfidfVectorizer` добавлены индикаторы утра, дня, вечера и ночи.
- Определены показатели метрики `roc_auc`.
- Обученные `LogisticRegressionCV` с найденными коэффициентами регуляризации применены для классификации сессий в соответствующих тестовых наборах данных.
- Отправлены результаты в соревновании на Kaggle.
- Сделаны выводы.

```
In (17): time_split = TimeSeriesSplit(n_splits=10) # для реализации
y_train = train_df['target'].astype('int').values
```

```
In (18): X_train_cnt_if = add_time_features(train_df.fillna(0), X_train)
X_test_cnt_if = add_time_features(test_df.fillna(0), X_test)
X_train_tfidf_if = add_time_features(train_df.fillna(0), X_train_tfidf)
X_test_tfidf_if = add_time_features(test_df.fillna(0), X_test_tfidf)
```

```
In (19): #time
logit_c_values = np.linspace(0.1, 5, 10)
logit_grid_searcher_cnt = LogisticRegressionCV(Cs=logit_c_values,
        solver='liblinear',
        random_state=17,
        cv_time_split=7,
        n_jobs=-1,
        scoring='roc_auc',
        max_iter=2000)
logit_grid_searcher_cnt.fit(X_train_cnt_if, y_train)
```

```
Out(19): [Parallel(n_jobs=-1)]: Using backend LokyBackend with 12 concurrent workers.
[Parallel(n_jobs=-1)]: Done 3 out of 10 | elapsed: 1.5min remaining: 1.5min
[Parallel(n_jobs=-1)]: Done 7 out of 10 | elapsed: 1.5min remaining: 1.5min
[Parallel(n_jobs=-1)]: Done 10 out of 10 | elapsed: 2.5min finished
```

```
In (20): logit_mean_cv_scores_cnt = logit_grid_searcher_cnt.scores_[1].mean(axis=0)
print(logit_mean_cv_scores_cnt.max())
logit_grid_searcher_cnt.get_params()
```

```
Out(20): 0.916699003305598 0.6444444444444444
logit_test_pred_cnt = logit_grid_searcher_cnt.predict_proba(X_test_cnt_if)[, 1]
write_to_submission_file(logit_test_pred_cnt, 'submit_cnt.csv') # 0.93969
```

```
In (22): #time
logit_c_values = np.linspace(0.1, 5, 10)
logit_grid_searcher_tfidf = LogisticRegressionCV(Cs=logit_c_values,
        solver='liblinear',
        random_state=17,
        cv_time_split=7,
        n_jobs=-1,
        scoring='roc_auc',
        max_iter=2000)
logit_grid_searcher_tfidf.fit(X_train_tfidf_if, y_train)
```

```
Out(22): [Parallel(n_jobs=-1)]: Using backend LokyBackend with 12 concurrent workers.
[Parallel(n_jobs=-1)]: Done 3 out of 10 | elapsed: 16.6s remaining: 38.8s
[Parallel(n_jobs=-1)]: Done 7 out of 10 | elapsed: 29.2s remaining: 12.3s
[Parallel(n_jobs=-1)]: Done 10 out of 10 | elapsed: 36.2s finished
```

```
Out(22): [Parallel(n_jobs=-1)]: Done 10 out of 10 | elapsed: 36.2s finished
logit_mean_cv_scores_tfidf = logit_grid_searcher_tfidf.scores_[1].mean(axis=0)
print(logit_mean_cv_scores_tfidf.max())
logit_grid_searcher_tfidf.get_params()
```

```
Out(23): 0.924549100379615 0.4555555555555556
logit_test_pred_tfidf = logit_grid_searcher_tfidf.predict_proba(X_test_tfidf_if)[, 1]
write_to_submission_file(logit_test_pred_tfidf, 'submit_tfidf.csv') # 0.94368
```

Поскольку более высокие `roc_auc` в leaderboard `TfidfVectorizer` соответствует более высокое значение данной метрики и на validation, и можно сделать вывод о неудачном выборе метода валидации. В дальнейшем будут использоваться наборы данных и преобразованные при помощи `TfidfVectorizer`.

## Применение новых фичей для классификации

```
In (25): X_train_if = X_train_cnt_if + X_train_tfidf_if
X_test_if = X_test_cnt_if + X_test_tfidf_if
```

```
In (26): X_train_2f = add_session_timespan(train_df, X_train_if)
X_test_2f = add_session_timespan(test_df, X_test_if)
```

```
In (27): X_train_3f = feature_short_session(train_df, X_train_2f)
X_test_3f = feature_short_session(test_df, X_test_2f)
```

```
In (28): X_train_4f = add_top_sites(train_df, X_train_3f)
X_test_4f = add_top_sites(test_df, X_test_3f)
```

```
In (29): X_train_5f = add_reverse_top_sites(train_df, X_train_4f)
X_test_5f = add_reverse_top_sites(test_df, X_test_4f)
```

```
In (30): X_train_6f = feature_middle_session(train_df, X_train_5f)
X_test_6f = feature_middle_session(test_df, X_test_5f)
```

```
In (44): # Провести cross-validation
logit_c_values = np.linspace(3, 5, 50)
logit_final = LogisticRegressionCV(Cs=logit_c_values,
        solver='liblinear',
        random_state=17,
        cv_time_split=7,
        n_jobs=-1,
        scoring='roc_auc',
        verbose=3,
        penalty='l2')
logit_final.fit(X_train_6f, y_train)
```

```
logit_mean_cv_scores = logit_final.scores_[1].mean(axis=0)
print(logit_mean_cv_scores.max())
logit_final.get_params()
```

```
Out(44): 0.9322137656775317 0.630612244897958 0.89044
[Parallel(n_jobs=-1)]: Using backend LokyBackend with 12 concurrent workers.
[Parallel(n_jobs=-1)]: Done 3 out of 10 | elapsed: 4.1min remaining: 9.6min
[Parallel(n_jobs=-1)]: Done 7 out of 10 | elapsed: 6.8min remaining: 2.9min
[Parallel(n_jobs=-1)]: Done 10 out of 10 | elapsed: 9.4min finished
```

```
In (45): logit_final = LogisticRegression(C=0.630612244897958,
        solver='liblinear',
        random_state=17,
        cv_time_split=7,
        n_jobs=-1,
        penalty='l2')
logit_final.fit(X_train_6f, y_train)
```

```
Out(45): LogisticRegression(C=0.630612244897958, n_jobs=-1, random_state=17,
        solver='liblinear')
logit_final_cv_score = cross_val_score(logit_final,
        X_train_6f,
        y_train,
        scoring='roc_auc',
        cv_time_split=7,
        verbose=3)
print('final cv score:', logit_final_cv_score)
```

```
In (47): logit_test_pred_final = logit_final.predict_proba(X_test_6f)[, 1]
write_to_submission_file(logit_test_pred_final, 'submit_final.csv')
```

## Заключение

Сам по себе анализ части тестовых данных: **0.95044**. Результат далек от лучших позиций в рейтинге, однако целью проделанной работы является знакомство с платформой Kaggle, получение навыков обработки и визуализации данных, а также изучение основ машинного обучения.

В качестве шагов для улучшения качества построенной модели в дальнейшем будут проделаны следующие шаги:

- применение других методов машинного обучения (бустинг над решающими деревьями, метод ближайших соседей);
- применение в схеме кросс-валидации (замечено, что не всегда рост `roc_auc` на кросс-валидации соответствует росту `roc_auc` на публичной части тестовых данных);
- проверка на переобучение модели при использовании отдельных признаков.

## Дополнительная глава (4)

В данном блоке собраны функции, позволяющие собирать данные с сайтов, которые пользователи используют. Сырые данные для этой задачи представляются как csv-файлы с данными о веб-сессиях отдельных пользователей в следующем виде:

```
timestamp,site
2013-11-15 11:40:34,google.com
...
```

Требуется объединить коллекцию таких файлов в одну таблицу данных. При этом в полученной таблице отдельные строки - это сессии - последовательности из нескольких сайтов. Решение должно поддерживать создание таблицы с разной длиной сессии, а также с различной шириной окна.

Пример: для длины сессии 10 и ширины окна 7 будет из 30 записей породить не 3 сессии (1-10, 11-20, 21-30), а 5 (1-10, 8-17, 15-24, 22-30, 29-30). При этом в предпоследней сессии будет один ноль, а в последней - 8 нулей.

Важное замечание: Задача, поставленная в этом разделе не связана напрямую с решением предложенной задачи соревнования на Kaggle. Задача решалась исключительно в академических целях и для того, чтобы осознать необходимость оптимизации кода.

```
In (1): import warnings
warnings.filterwarnings('ignore')
from glob import glob

import os
import pickle
from tqdm import tqdm
import numpy as np
import pandas as pd

import itertools

import re
import scipy.sparse as sp

import time

PATH_TO_DATA = os.path.join('initial_data', 'users')
PATH_TO_DICT_FREQ = os.path.join('initial_data', 'site_freq')
```

Посмотрим на файл с данными на примере пользователя user0128 из коллекции с 10 пользователями.

```
In (2): user31_data = pd.read_csv(os.path.join(PATH_TO_DATA, '10users/user0128.csv'))
user31_data.head()
```

```
Out(2): timestamp site
0 2013-11-15 13:4603 fpdfdownload2.mactrom.com
1 2013-11-15 13:4613 mail.google.com
2 2013-11-15 13:4613 www.gmail.com
3 2013-11-15 13:4625 accounts.google.com
4 2013-11-15 13:4628 accounts.youtube.com
```

Основная функция в этом разделе - это функция

```
prepare_sparse_train_set_window.
```

Она принимает на вход путь к коллекции csv-файлов с данными по каждому пользователю, путь к заранее подготовленному словарю сайтов, а также параметры `session_length` и `window_size`, отвечающие за длину сессии и размер окна соответственно.

Для удобства, отдельно определена функция



```
In [4]: def prepare_data_set_window(path_to_csv_files, site_freq_path,
                                session_length=10, window_size=10):
    """Подготавливаем набор данных. Самих данных хранятся в path_to_csv_files.
    Возвращает матрицу частот в разреженном формате, словарь с user_id и DataFrame
    соответствующих наборов сессий"""

    with open(site_freq_path, 'rb') as file:
        site_freq_dict = pickle.load(file)

    data = []

    list_of_files = glob(os.path.join(path_to_csv_files, ''*))
    for path_to_user in list_of_files:
        user_id = int(re.sub(r'\b\d+', '', path_to_user[-8:-4], 1)) # user0001 --> 1
        sites_array = pd.read_csv(path_to_user)['site'].map(lambda x: site_freq_dict[x][0]).values.tolist()
        n = len(sites_array)
        ind = 0
        while True:
            if ind + session_length > n-1:
                data.append(sites_array[ind:n] + [0 for _ in range(ind + session_length - n)] + [user_id])
            else:
                data.append(sites_array[ind:ind+session_length] + [user_id])
            ind += window_size
            if ind >= n:
                break

    feature_names=['site%i'%i for i in range(1, session_length + 1)] + ['target']
    data_df = pd.DataFrame(data, columns=feature_names)
    target = np.array(data_df['target']).values, dtype='int16')

    """Создание и заполнение разреженной матрицы частот"""
    X, y = data_df.iloc[:, :-1].values, data_df.iloc[:, -1].values
    sparse_data = sp.csr_matrix(make_sparse_data(X), dtype='int8')

    return sparse_data, target, data_df
```

```
In [5]: %time
X_sparse_10users_s10_w7, y_10users_s10_w7, df = prepare_data_set_window(os.path.join(PATH_TO_DATA, '150users'),
                                os.path.join(PATH_TO_SITE_FREQ, 'site_freq_150users.pkl'),
                                session_length=10, window_size=7)
```

```
100% | 195712/195712 | 00:06<00:00,
30078.81it/s |
Wall time: 9.7 s

Обрабатываем набор данных для 10 и 150 пользователей и сохраням разреженные матрицы в формате
```

```
X_sparse_(кол-во пользователей)users_(длина сессии)_w(ширина окна).pkl

Также сохраням id пользователей в формате
```

```
y_{кол-во пользователей}users_(длина сессии)_w(ширина окна).pkl

In [6]: %time
for num_users in [10, 150]:
    for window_size, session_length in itertools.product([10, 7, 5], [15, 10, 7, 5]):
        if window_size <= session_length:
            X_sparse, y, _ = prepare_data_set_window(os.path.join(PATH_TO_DATA, '{users}'.format(num_users)),
                                                    os.path.join(PATH_TO_SITE_FREQ, 'site_freq_{users}_{
                                                    session_length}session_length_{window_size}window_size.pkl'),
                                                    'X_sparse_{users}_s{w}.pkl'.format(num_users,
                                                    session_length,
                                                    window_size)), 'wb') as X_pkl:
                pickle.dump(X_sparse, X_pkl)
            with open(os.path.join(PATH_TO_DATA,
                                    'y_{users}_s{w}.pkl'.format(num_users,
                                    session_length,
                                    window_size)), 'wb') as y_pkl:
                pickle.dump(y, y_pkl)
```

```
100% | 14061/14061 | 00:08<00:00,
29970.89it/s |
100% | 14061/14061 | 00:08<00:00,
30358.86it/s |
100% | 20087/20087 | 00:08<00:00,
30059.95it/s |
100% | 20087/20087 | 00:08<00:00,
30959.09it/s |
100% | 20087/20087 | 00:08<00:00,
31821.81it/s |
100% | 28118/28118 | 00:09<00:00,
29757.74it/s |
100% | 28118/28118 | 00:09<00:00,
30645.10it/s |
100% | 28118/28118 | 00:08<00:00,
31796.35it/s |
100% | 28118/28118 | 00:08<00:00,
32044.28it/s |
100% | 137019/137019 | 00:04<00:00,
28481.56it/s |
100% | 137019/137019 | 00:04<00:00,
29834.93it/s |
100% | 195712/195712 | 00:06<00:00,
29249.86it/s |
100% | 195712/195712 | 00:06<00:00,
30449.86it/s |
100% | 195712/195712 | 00:06<00:00,
30198.01it/s |
100% | 273957/273957 | 00:09<00:00,
29314.61it/s |
100% | 273957/273957 | 00:08<00:00,
32020.28it/s |
100% | 273957/273957 | 00:08<00:00,
31819.93it/s |
100% | 273957/273957 | 00:08<00:00,
32502.36it/s |
Wall time: 1min 39s
```

Теперь, когда подготовлены таблицы с сессиями в различных форматах (с различными session\_length и window\_size), можно попоробовать натренировать модели (например, логистическую регрессию) на различных наборах данных. Тренировать модели на данных 150 пользователей долго, поэтому для примера, воспользуемся датасетами для 10 пользователей.

```
In [7]: from sklearn.linear_model import LogisticRegressionCV
from sklearn.model_selection import StratifiedKFold

def logit_cv(logit_c_values, path_to_X_pickle, path_to_y_pickle):
    """ Тренируем LogisticRegressionCV с параметром на logit_c_values
    на данных из path_to_X_pickle и path_to_y_pickle"""

    with open(path_to_X_pickle, 'rb') as X_sparse_users_pkl:
        X_sparse_users = pickle.load(X_sparse_users_pkl)
    with open(path_to_y_pickle, 'rb') as y_pkl:
        y = pickle.load(y_pkl)

    logit_c_values = np.linspace(0.5, 4, 5)

    logit_grid_searcher = LogisticRegressionCV(Cs=logit_c_values,
                                                multi_class='multinomial',
                                                random_state=17,
                                                cv=StratifiedFold(n_splits=3, shuffle=True, random_state=17),
                                                n_jobs=-1)

    logit_grid_searcher.fit(X_sparse_users, y)

    return logit_grid_searcher

logit_c_values = np.linspace(0.5, 6, 5)
```

```
In [8]: all_logit_searchers = []

for window_size, session_length in itertools.product([10, 7, 5], [15, 10, 7, 5]):
    if window_size <= session_length:
        path_to_X_pickle = os.path.join(PATH_TO_DATA,
                                         'X_sparse_{users}_s{session_length}_w{window_size}.pkl')
        path_to_y_pickle = os.path.join(PATH_TO_DATA,
                                         'y_{users}_s{session_length}_w{window_size}.pkl')

        logit_grid_searcher = logit_cv(logit_c_values, path_to_X_pickle, path_to_y_pickle)

        all_logit_searchers.append(logit_grid_searcher)

    print(f'session_length = {session_length}, window_size = {window_size}')
    logit_mean_cv_scores = logit_grid_searcher.scores_[3].mean(axis=0)
    print(f'max cv score for user0031: ', logit_mean_cv_scores.max())
    print(f'best C: ', logit_grid_searcher.Cs_[logit_mean_cv_scores.argmax()])
```

```
session_length = 15, window_size = 10
max cv score for user0031:  0.8381338453879525
best C:  1.375
session_length = 10, window_size = 10
max cv score for user0031:  0.7763316976032999
best C:  1.375
session_length = 15, window_size = 7
max cv score for user0031:  0.8637924201474676
best C:  3.125
session_length = 10, window_size = 7
max cv score for user0031:  0.81320162268245263
best C:  1.375
session_length = 7, window_size = 7
max cv score for user0031:  0.758749404040411881
best C:  3.125
session_length = 15, window_size = 5
max cv score for user0031:  0.8841665915712135
best C:  3.125
session_length = 10, window_size = 5
max cv score for user0031:  0.8334162722580608
best C:  4.0
session_length = 7, window_size = 5
max cv score for user0031:  0.7878939587540709
best C:  3.125
session_length = 5, window_size = 5
max cv score for user0031:  0.7368590432312572
best C:  3.125
```

Наилучшее качество на кросс-валидации оказалось у датасета session\_length = 15, window\_size = 5. В дальнейшем можно проверить то же самое на более крупных датасетах, например на датасете 3000users. Но это выходит за рамки дополнительной главы, которая на этом заканчивается.