

Оценка зависимости от социальной сети VK

Цель исследования: по открытым данным оценить зависимость подростка-студента от социальной сети VK.

Источник данных: страница подростка в VK.

Получение данных: web-scraper загружает страницу каждые 15 минут и получает статус пользователя (online, offline, заходил x мин. назад и т.д.)

Этика: слежка за подростком неэтична и незаконна. Прибегать к ней возможно тогда и только тогда, когда родителей беспокоит здоровье или безопасность их ребенка.

Законодательство: согласно статье 152.2 ГК РФ и статье 137 УК РФ сбор, хранение и распространение информации о частной жизни гражданина без его согласия запрещены и карается штрафом или лишением свободы. Уровень доступности страницы VK для публичного просмотра, установленный пользователем в настройках безопасности ("весь интернет", "только пользователи VK", "только друзья") можно рассматривать не более как согласие пользователя на получение информации о нём (соответствующим кругом лиц), но ни в коем случае не как разрешение хранить или распространять её.

В данном исследовании используются анонимизированные данные случайного пользователя VK.

План исследования:

- 1. Подготовка данных.** Данные будут загружены из лог-файла скрапера, затем проверены и преобразованы в формат, удобный для исследования.
- 2. Анализ.** Будет выполнен анализ данных с целью установить наличие, степень и динамику зависимости от социальной сети.
- 3. Итоги.** На основе анализа будут сделаны окончательные выводы о зависимости от социальной сети.

Подготовка данных

Загрузка

Подключим необходимые библиотеки и загрузим данные из лог-файла скрапера:

```
In [1]: import datetime
import numpy as np
import pandas as pd
import matplotlib.pyplot as plt

time_unit = 15 # 15 минут - единица временной шкалы
day_name_ru = ('пн','вт','ср','чт','пт','сб','вс')
month_name_ru = ('---', 'январь', 'февраль', 'март', 'апрель', 'май',
                 'июнь', 'июль', 'август', 'сентябрь', 'октябрь', 'ноябрь', 'декабрь')

log = pd.read_csv('portfolio-log.txt', sep='\\t', header=None)
log.columns = ['datetime', 'is_online', 'min_left', 'comment']

Проверим, корректно ли загружены данные:
```

```
In [2]: log.head(5)

Out[2]:
```

	datetime	is_online	min_left	comment
0	03.08.2021 08:00	NaN	NaN	[m] заходил сегодня в 1:57
1	03.08.2021 08:15	NaN	NaN	[m] заходил сегодня в 1:57
2	03.08.2021 08:30	NaN	NaN	[m] заходил сегодня в 1:57
3	03.08.2021 08:45	NaN	NaN	[m] заходил сегодня в 1:57
4	03.08.2021 09:00	NaN	NaN	[m] заходил сегодня в 1:57

```
In [3]: log.shape

Out[3]: (2480, 4)
```

Структура и значения в целом соответствуют ожидаемым.

Итого: данные загружены успешно.

Предобработка

Исправим типы данных в столбцах:

```
In [4]: log['datetime'] = pd.to_datetime(log['datetime'], format='%d.%m.%Y %H:%M')
log['is_online'] = log['is_online'].astype('int64')
log['min_left'] = log['min_left'].astype('int64')
log.dtypes

Out[4]:
```

	datetime	is_online	min_left	comment
	datetime64[ns]	bool	int64	object
	min_left	bool	int64	object
	comment	object		
	dtype:	object		

Проверим результат изменения типов данных:

```
In [5]: log.tail(3)

Out[5]:
```

	datetime	is_online	min_left	comment
2477	2021-09-11 19:15:00	False	<NA>	заходил час назад
2478	2021-09-11 19:30:00	True	0	NaN
2479	2021-09-11 19:45:00	False	15	NaN

Нормализуем текст комментариев и удалим записи, где скрапер сообщает о возникновении сбоев (и невозможности получить данные):

```
In [6]: # Нормализуем текст комментариев
log['comment'] = log['comment'].str.strip().str.lower()
# Отираем номера строк со сбоями загрузки
exception_rows_mask = log[log['comment'].str.startswith('exception', na=False)].index
# Удалим строки со сбоями
log = log.drop(exception_rows_index)

Данные лог-файла скрапера приведены к нужным типам, из них исключены записи со сбоями.
```

Итого: данные лог-файла скрапера готовы к использованию.

Реструктуризация

Заранее известно, что лог-файл скрапера содержит пропуски. Во-первых, в период с 2:00 до 8:00 скрапер не собирает данные. Во-вторых, скрапер временно может быть отключен или не запущен по ошибке. Соответственно в лог-файле не окажется записей за то время, когда скрапер не работал. Для анализа данных это неудобно. Реструктурируем данные так, чтобы на каждые 15 минут каждые сутки было по одной записи в наборе данных:

```
In [7]: min_date = log['datetime'].min().floor('D') # 22.01.2021 11:00 --> 22.01.2021 00:00
max_date = log['datetime'].max().floor('D') # 30.01.2021 18:45 --> 30.01.2021 00:00
data = []
dt = min_date
while dt < max_date:
    data.append(dt)
    dt += datetime.timedelta(minutes=15)

# Создадим набор данных без пропусков времени и дат (одна строка на time_unit минут)
df = pd.DataFrame(index=data, columns=['datetime'])
# Добавим данные из лога (на те моменты времени, которые есть в логе)
df = df.merge(log, on='datetime', how='left')
# Пометим на будущее в отдельной колонке, какие данные есть, каких нет
df['is_na'] = df['is_online'].isna()

Проверим временную структуру данных:
```

```
In [8]: df.shape[0], df['datetime'].min(), df['datetime'].max()

Out[8]: (3748, Timestamp('2021-08-03 00:00:00'), Timestamp('2021-09-10 23:45:00'))
```

Итого: временная структура данных исправлена.

Пропуски

Восстановим пропуски по данным, о том сколько минут назад был пользователь в сети:

```
In [9]: x = (
df[df['min_left'] > time_unit]
df.apply(lambda row: row['datetime'] - datetime.timedelta(minutes=row['min_left']), axis=1)
    .dt.ceil(freq=str(time_unit) + 'T')
    .unique()
)

rows_mask = df['datetime'].isin(x) & df['is_online'].isna()
df.loc[rows_mask, ['is_online', 'is_na']] = [True, False]
filled_by_minutes = df[rows_mask].index.to_list() # сохраняем для справки

Выделим компоненты даты и времени в отдельные столбцы для удобства дальнейшей обработки:
```

```
In [10]: df['y'] = df['datetime'].dt.year
df['m'] = df['datetime'].dt.month
df['d'] = df['datetime'].dt.day
df['h'] = df['datetime'].dt.hour
df['W'] = df['datetime'].dt.minute

По фразе "заходил сегодня в чч:мм" восстановим пропуски в данных (где это возможно и нужно):
```

```
In [11]: x = (
df[df['comment'].notna() & df['comment'].str.contains('заходил сегодня')]
    .pivot_table(index=['y','m','d','comment'], values='datetime', aggfunc='first')
    .reset_index()
    .apply(lambda row: row['datetime'].row['y'], row['m'], row['d'],
        int(row['comment'][-5:-3]), int(row['comment'][-2:-1]))
    .dt.ceil(freq=str(time_unit) + 'T')
    .unique()
)

rows_mask = df['datetime'].isin(x) & df['is_online'].isna()
df.loc[rows_mask, ['is_online', 'is_na']] = [True, False]
filled_by_comments = df[rows_mask].index.to_list() # сохраняем для справки

Заполним оставшиеся пропуски столбца is_online значениями False:
```

```
In [12]: # Заполняем ночные пропуски, когда сканирование не работает
rows_mask = (df['H'] >= 2) & (df['H'] <= 7) & df['is_online'].isna()
filled_by_nights = df[rows_mask].index.to_list()
df.loc[rows_mask, ['is_online', 'is_na']] = [False, False]
# Заполняем остальные пропуски
df['is_online'] = df['is_online'].fillna(False).astype(bool)

Удалим ненужные столбцы:
```

```
In [13]: df = df.drop(columns=['min_left', 'comment'])

Проверим наличие пропусков:
```

```
In [14]: df.isna().sum()

Out[14]:
```

	datetime	is_online	is_na	y	m	d	h	m	dtype
	0	0	0	0	0	0	0	0	int64

Итого: пропуски заполнены, пропусков нет.

Набор данных с и без пропуска дней

Вероятны ситуации, когда в какие-то дни данных о посещениях социальной сети нет вообще (скрапер целые сутки не работал) или записей слишком мало (скрапер работал не полные сутки или долгое время сеть была недоступна). Условимся, что если пропущено более 20% записей за сутки (около 5 часов), то такие сутки следует считать пропущенными (для большинства видов анализа). В дальнейшем удобно иметь два набора данных: с пропущенными днями и без. Создадим такие наборы:

```
In [15]: df_all_days = df.copy()
maxnap = 0.28 # макс. число пропущенных записей, %
na_col = 'N/A'
gindex = ['y','m','d']
na = df.groupby(gindex)['is_na'].sum().reset_index().rename(columns={'is_na':na_col})
na[na_col] = na[na_col] / (24 * 60 / time_unit)
df = df.merge(na, on=gindex, how='left')
df = df[df[na_col] <= maxnap]
df = df.drop(columns=[na_col])
days_deleted = na[na[na_col] > maxnap]
df_filtred_by_na = df.copy()

In [16]: print(f'Удалено дней с критичным пропуском данных (более {0.0%}) - {0}', format(maxnap, days_deleted.shape[0]))
print()
print(f'Итого: df_all_days - набор со всеми днями независимо от числа пропусков, {df_all_days.shape[0]} строк')
print(f'Итого: df_filtred_by_na - набор без дней с большими пропусками, {df_filtred_by_na.shape[0]} строк')
print()
```

Удалено дней с критичным пропуском данных (более 20%) - 7.

df_all_days - набор со всеми днями независимо от числа пропусков, 3748 строк
df_filtred_by_na - набор без дней с большими пропусками, 3676 строк

Итого: наборы с пропуском дней и без созданы.

ИТОГИ

Подведём итоги подготовки данных:

```
In [17]: df = df_filtred_by_na.copy()
begin = df['datetime'].min()
end = df['datetime'].max()
period_len_days = df.groupby(['y','m','d']).ngroups
fake_offline_n = df['is_na'].sum()
print(
    '---\n'
    'Период наблюдений: {} - {}'.format(begin, end)
    '---\n'
    'Дней с числом пропусков в пределах допустимого - {}'.format(len(days_deleted))
    'Записей в логе - {}, в наборе - {}'.format(len(log), len(df))
    'Заполнено пропусков (по min_left/no comment/ночь) - {}/{}(/{})\n'.format(len(log), len(df))
    'Записей fake offline (для них возможно, что реальный статус online) - {}'.format(len(days_deleted))
    '---\n'
    'Итого: df_all_days - набор со всеми днями независимо от числа пропусков, {df_all_days.shape[0]} строк'
    'Итого: df_filtred_by_na - набор без дней с большими пропусками, {df_filtred_by_na.shape[0]} строк'
    '---\n'
    'Итого: наборы с пропуском дней и без созданы.'
)

Создадим подпрограмму, вычисляющую распределение времени суток на online, offline, условный-сон для заданного множества дней:
```

```
In [18]: def get_sleep8(df):
    gindex = ['y','m','d']
    x = df[df['is_online'] & (df['H'] < 5)].pivot_table(
        index=gindex, values='datetime', aggfunc='max')
    y = df[df['is_online'] & (df['H'] >= 5) & (df['H'] < 8)].pivot_table(
        index=gindex, values='datetime', aggfunc='min')
    m = df[gindex].groupby(gindex).first().reset_index()
    if not x.empty:
        m = m.merge(x, on=gindex, how='left').rename(columns={'datetime':'tx'})
    else:
        m['tx'] = np.datetime64('NaT')
    m.loc[m['tx'].isna(), 'tx'] = m['tx'].isna().apply(axis=1, func=lambda row:
        datetime.datetime(row['y'], row['m'], row['d'], 8, 0, 0))
    if not y.empty:
        m = m.merge(y, on=gindex, how='left').rename(columns={'datetime':'ty'})
    else:
        m['ty'] = np.datetime64('NaT')
    m.loc[m['ty'].isna(), 'ty'] = m['ty'].isna().apply(axis=1, func=lambda row:
        datetime.datetime(row['y'], row['m'], row['d'], 0, 0, 0))
    m['ty - tx'] = (m['ty'] - m['tx']).dt.total_seconds()
    #print(m)
    sum_ticks = int(m['ty - tx'].sum() / 60 / time_unit)
    ave_ticks = sum_ticks / m.shape[0]
    return (sum_ticks, ave_ticks)

Создадим подпрограмму, вычисляющую распределение времени суток на online, offline, условный-сон для заданного множества дней:
```

```
In [19]: def get_pie_parts(df, start, end):
    data = df[(df['datetime'] > start) & (df['datetime'] <= end)]
    days = data.groupby(['y','m','d']).ngroups
    online_mean = data['is_online'].sum() / days * time_unit / 60
    sleep8_mean = get_sleep8(data)
    sleep8_mean *= time_unit / 60
    offline924_mean = 24 - online_mean - sleep8_mean;
    online_mean = round(online_mean,1)
    sleep8_mean = round(sleep8_mean,1)
    offline924_mean = round(offline924_mean,1)
    return (online_mean, offline924_mean, sleep8_mean, days)

Получим среднее распределение времени в сутках для всего периода наблюдений, за последние 30 дней и за последние 7 дней:
```

```
In [20]: df = df_filtred_by_na.copy()
begin = df['datetime'].min()
end = df['datetime'].max()
pie_30_days = get_pie_parts(df, begin, end)
pie_7_days = get_pie_parts(df, end - datetime.timedelta(days=30), end)
pies_df = pd.DataFrame([pie_all_days, pie_30_days, pie_7_days],
    columns=['online', 'offline', 'sleep8', 'days rec'],
    index = ['all', '30d', '7d'])
pies_df.index.name = 'period (last n days)'
pies_df
```

```
Out[20]:
```

	online	offline	sleep8	days rec
period (last n days)				
all	10.6	6.6	6.9	32
30	10.9	6.1	7.0	24
7	11.9	4.8	7.4	7

Отобразим таблицу в графическом виде:

```
In [21]: pies_df = pies_df.drop(columns=['days rec'])
pies_df.plot(subplots=True, figsize=(12, 4),
    legend=False, explode=[0.1,0.0], autopct='%1.0f%',
    colors=['salmon', 'skyblue', 'steelblue'])
title = "Распределение времени в сутках (за всё время, за последние 30дней и последние 7 дней)"
plt.show()
```

Распределение времени в сутках (за всё время, за последние 30дней и последние 7 дней)

Выводы:

1. Почти половину суток составляют 15-минутные отрезки, когда подросток заходил в сеть не менее одного раза (особенно ярко это проявляется в последние семь дней).
2. Времена условного сна (00:00-08:00) вместо 8 часов составляет в среднем 7 часов за период наблюдений и лишь с началом учебного года стало наблюдаться меньше чем на 15 минут).
3. 15-минутные отрезки времени с 8:00 до 00:00, когда подросток не был.сет, в сумме за период составляют лишь 5-6 часов - почти в два раза меньше, чем время online.

Статистика по дням

Рассмотрим, как изменялось время присутствия в сети на протяжении дней:

```
In [22]: df = df_filtred_by_na.copy()
# Отмечаем начало выхода в онлайн (чтобы потом получить число сессий за сутки)
df['online_start'] = (df['is_online'] - df['is_online'].shift()) > 0
gindex = ['y','m','d']
days = (
    df
    .pivot_table(index=gindex, values=['is_online', 'is_na', 'online_start'], aggfunc='sum')
    .reset_index()
)
days['date'] = days.apply(axis=1, func=lambda row: datetime.date(row['y'], row['m'], row['d']))
days['online'] = days['is_online'] * time_unit / 60
days['offline'] = 24 - days['online']
days['is_na'] = days['is_na'] * time_unit / 60
days['rolling_mean'] = days['online'].rolling(7, min_periods=7).mean()
# Средняя длительность сессии онлайн (на протяжении более чем на time_unit минут)
days['online_mean'] = days['online'] / days['online_start']

cardio = days[['date', 'is_na', 'online_mean', 'offline_mean', 'rolling_mean']]
ax = cardio.plot(figsize=(min(25, days.shape[0]*0.25), 5), stacked=False,
    rot=20, legend=True, color=['#E6E6FA', '#FFA07A', '#FF69B4', '#FF69B4', '#FF69B4'])
ax.grid(True, style='dotted')
ax.set_title('История присутствия в социальной сети')
ax.set_ylim(-1.0, 10.0)
ax.set_xticks(ticks=[0,1,2,4,6,7,8,10,12,13,14,16,18])
ax.set_ylabel('часов')
ax.set_xlabel('')
ax.legend(['пропущено данных', 'средний сеанс online', 'суммарное время online',
    'скользящее среднее за 7-дней'],
    loc='upper left')
plt.show()
```

История присутствия в социальной сети

Выводы:

1. История по дням наблюдается привычка регулярно засиживаться в сети за полночь, нередко позже часа ночи.
2. Первые учебные дни не привели к тому, что подросток на время занятий способен отвлечься от социальной сети. Периоды с режими дневания в сеть появляются лишь на второй учебной неделе и то лишь во второй половине дня или вечером и не во все учебные дни.

Статистика по дням недели

Посмотрим как распределено время присутствия в сети по дням недели:

```
In [29]: days = (
df_filtred_by_na[['datetime', 'is_online']]
    .copy()
    .set_index('datetime')
    .resample('1D') # включает в days дни, которых нет в df_filtred_by_na
    .agg('sum')
)
days = days.rename(columns={'is_online':'online'})
days['date'] = days.drop(days.index[days['online'] == 0].index) # исправляем побочный эффект .resample('1D')
days['weekday'] = days.index.weekday
days['online'] = days['online'] * 15 / 60
days['offline'] = 24 - days['online']
weekdays = days.pivot_table(index='weekday', values=['online', 'offline'], aggfunc='mean')
weekdays = weekdays[['online', 'offline']]
weekdays['day_name'] = weekdays.apply(axis=1, func=lambda row: day_name_ru[row.name])

fig, axes = plt.subplots(nrows=4, ncols=2, figsize=(12, 3),
    values='is_online', aggfunc='sum')

ax = weekdays28.plot(ax=axes[1], width=0.6,
    kind='bar', stacked=True, x='day_name', xlabel='', rot=0,
    legend=True, color=['salmon', 'skyblue'])
ax.set_title('Присутствие в VK по дням недели\n(за последние 28 дней)')
ax.grid(axis='y')
ax.set_ylim(0, 24)
ax.set_xticks(ticks=[4,8,12,16,20])
ax.legend(['upper left', 'bbox_to_anchor=(1.0, 1.0)'])
ax.set_ylabel('часов')

ax = weekdays.plot(ax=axes[0], width=0.7,
    kind='bar', stacked=True, x='day_name', xlabel='', rot=0,
    legend=True, color=['salmon', 'skyblue'])
ax.set_title('Присутствие в VK по дням недели\n(за весь период наблюдений)')
ax.grid(axis='y')
ax.set_ylim(0, 24)
ax.set_xticks(ticks=[4,8,12,16,20])
ax.legend(['upper left', 'bbox_to_anchor=(1.0, 1.0)'])
ax.set_ylabel('часов')

plt.show()
```

Присутствие в VK по дням недели (за весь период наблюдений)

Присутствие в VK по дням недели (за последние 28 дней)

Выводы. Большое различие между днями недели не наблюдается. Нельзя сказать, что по поведению подростка можно отличить, например, будние дни от выходных. Присутствие в VK во все дни недели примерно одинаково. Распределение по дням недели не позволяет сделать никаких определённых выводов.

Итоги исследования

Данные показывают, что у подростка присутствует сильная зависимость от социальной сети VK. Суммарно за сутки он проводит в ней (не отвлекаясь более, чем на 15 минут) по 10-12 часов, регулярно ложится спать глубоко заспанным и постоянно отвлекается от занятий или вовсе не посещает их. По реакции на начало учебного года (аномально длительное время присутствия в сети в течение первой недели) есть основания полагать, что подросток живому общению предпочитает жизнь в сети.

In []: