

In [2]:

```
1 import matplotlib.pyplot as plt
2 import numpy as np
3 import pandas as pd
4 import plotly.express as px
5 import scipy.stats as sps
6 import seaborn as sns
7
8
9 sns.set(font_scale=1.5)
10 %matplotlib inline
```

executed in 5.96s, finished 13:19:40 2020-05-22

**Формулировка задачи на простом "пользовательском" языке:\**

Мы займемся анализом поездок пользователей в зависимости от типа абонента/подписки. Как используют велопрокат пользователи типа Member и Short-Term Pass Holder .\

**Формулировка задачи на математическом языке:\**

Выяснить, как тип пользователя влияет на продолжительность поездки. Узнать, есть ли различие между средним значением продолжительности поездки для подписки типа Member и Short-Term Pass Holder .

Загрузим чуть-чуть обработанные данные, где удалены дубликаты первых 50000 поездок в таблице trips.csv.

In [3]:

```
1 DATA_DIR = '../data/processed/'
```

executed in 9ms, finished 13:19:40 2020-05-22

In [4]:

```
1 stations_data = pd.read_csv(DATA_DIR+"station.csv")
2 trips = pd.read_csv(DATA_DIR+"trips.csv", error_bad_lines=False, index_col=0)
3 weather = pd.read_csv(DATA_DIR+"weather.csv")
```

executed in 2.42s, finished 13:19:43 2020-05-22

In [5]:

```
1 trips.head(3)
```

executed in 101ms, finished 13:19:43 2020-05-22

Out[5]:

	trip_id	starttime	stoptime	bikeid	tripduration	from_station_name	to_station_name
0	431	10/13/2014 10:31	10/13/2014 10:48	SEA00298	985.935	2nd Ave & Spring St	Occident S
1	432	10/13/2014 10:32	10/13/2014 10:48	SEA00195	926.375	2nd Ave & Spring St	Occident S
2	433	10/13/2014 10:33	10/13/2014 10:48	SEA00486	883.831	2nd Ave & Spring St	Occident S

Возможны два типа подписки.

In [6]:

```
1 trips.usertype.unique()
```

executed in 77ms, finished 13:19:44 2020-05-22

Out[6]:

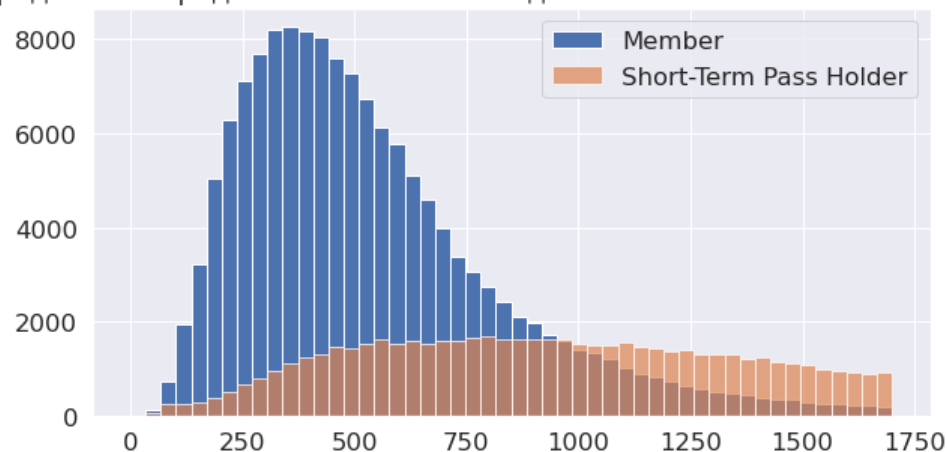
array(['Member', 'Short-Term Pass Holder'], dtype=object)

In [7]:

```
1 plt.figure(figsize=(10, 5))
2 plt.title('Распределение продолжительности поездки в зависимости от типа абонента')
3 plt.hist(trips[trips['usertype'] == 'Member'].tripduration,
4          bins=50, range=(0, 1700),
5          label='Member')
6 plt.hist(trips[trips['usertype'] == 'Short-Term Pass Holder'].tripduration,
7          bins=50, range=(0, 1700),
8          label='Short-Term Pass Holder', alpha=0.7)
9 plt.legend()
10 plt.show()
```

executed in 2.39s, finished 13:19:46 2020-05-22

Распределение продолжительности поездки в зависимости от типа абонента



Посмотрим на статистики для tripduration для разных типов

In [8]:

```
1 trips.drop(['trip_id', 'birthyear'], axis=1).groupby(by=['usertype']).describe()
```

executed in 486ms, finished 13:19:47 2020-05-22

Out[8]:

tripduration								
	count	mean	std	min	25%	50%	75%	max
usertype								
Member	146171.0	592.977313	731.550556	60.008	324.2025	479.1490	693.39150	27985.8
Short-Term Pass Holder	89894.0	2193.901312	2959.113053	60.111	760.0590	1251.4325	2161.54875	28794.3

Можем проверить критерием Уилкоксона-Манна-Уитни гипотезу об однородности.

In [9]:

```
1  sps.mannwhitneyu(trips[trips['usertype'] == 'Member'].tripduration,  
2                  trips[trips['usertype'] == 'Short-Term Pass Holder'].tripduration)
```

executed in 437ms, finished 13:19:48 2020-05-22

Out[9]:

MannwhitneyUResult(statistic=2077197774.5, pvalue=0.0)

Теперь посмотрим, когда разные типы пользователей пользуются велосипедом. Построим гистограммы времени начала поездки.

In [11]:

```
1  def is_weekend(day): ↵
```

executed in 12ms, finished 13:21:32 2020-05-22

In [15]:

```
1  trips['starttime'] = pd.to_datetime(trips['starttime'])  
2  trips['stoptime'] = pd.to_datetime(trips['stoptime'])  
3  trips['starttime_of_day_hour'] = [hour for hour in trips['starttime'].dt.hour]  
4  trips['weekend'] = is_weekend(trips.starttime.dt.weekday)
```

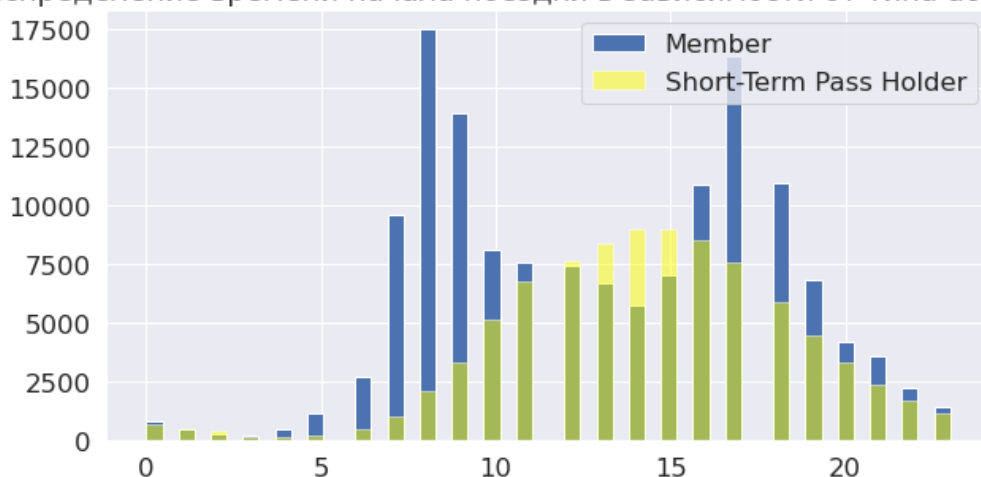
executed in 287ms, finished 13:22:13 2020-05-22

In [13]:

```
1  plt.figure(figsize=(10, 5))  
2  plt.title('Распределение времени начала поездки в зависимости от типа абонемента')  
3  plt.hist(trips[trips['usertype'] == 'Member'].starttime_of_day_hour,  
4          bins=50,  
5          label='Member')  
6  plt.hist(trips[trips['usertype'] == 'Short-Term Pass Holder'].starttime_of_day_hour,  
7          bins=50,  
8          label='Short-Term Pass Holder', alpha=0.5, color='yellow')  
9  plt.legend()  
10 plt.show()
```

executed in 889ms, finished 13:21:39 2020-05-22

Распределение времени начала поездки в зависимости от типа абонемента

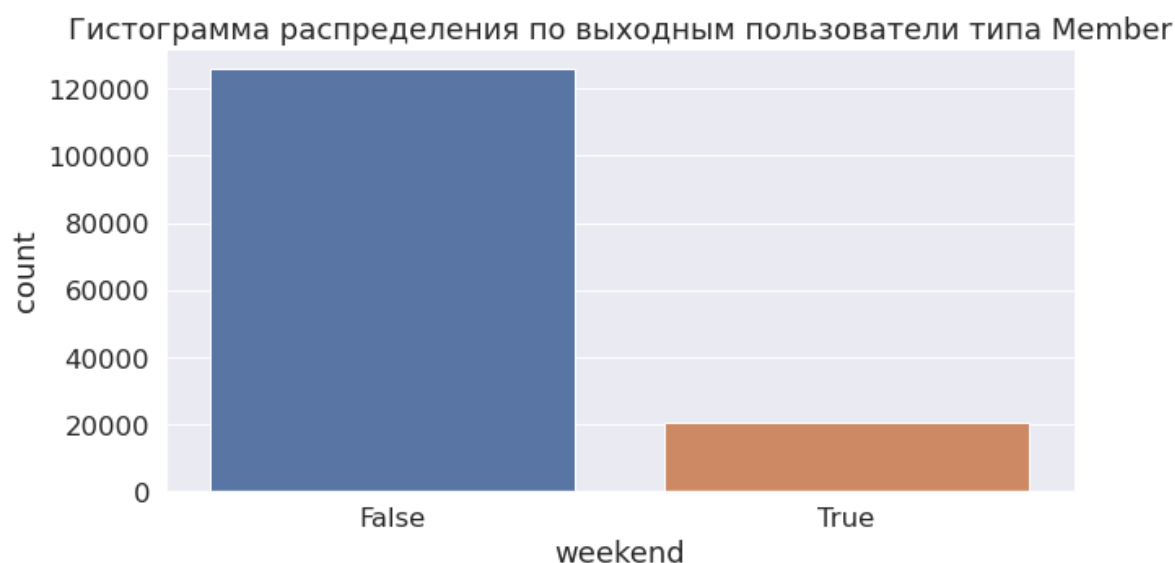


Посмотрим, когда чаще всего катаются Member , Short-Term Pass Holder .

In [18]:

```
1 plt.figure(figsize=(10, 5))
2 plt.title('Гистограмма распределения поездок по выходным пользователи типа Member')
3 sns.countplot(x='weekend', data=trips[trips['usertype'] == 'Member'])
4 plt.show()
```

executed in 273ms, finished 13:23:14 2020-05-22

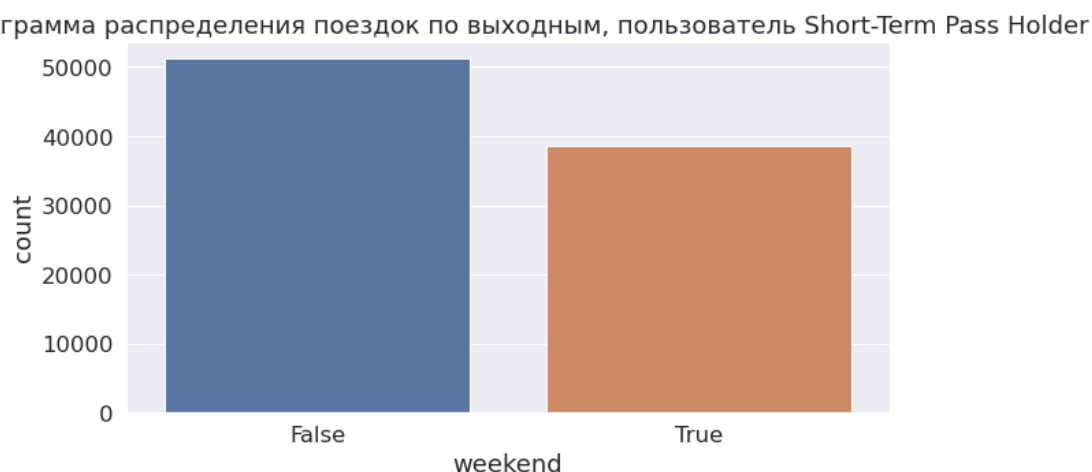


Как видим, Member чаще катаются в будний день.

In [20]:

```
1 plt.figure(figsize=(10, 5))
2 plt.title('Гистограмма распределения поездок по выходным, пользователь Short-Term Pass Holder')
3 sns.countplot(x='weekend', data=trips[trips['usertype'] != 'Member'])
4 plt.show()
```

executed in 260ms, finished 13:24:25 2020-05-22



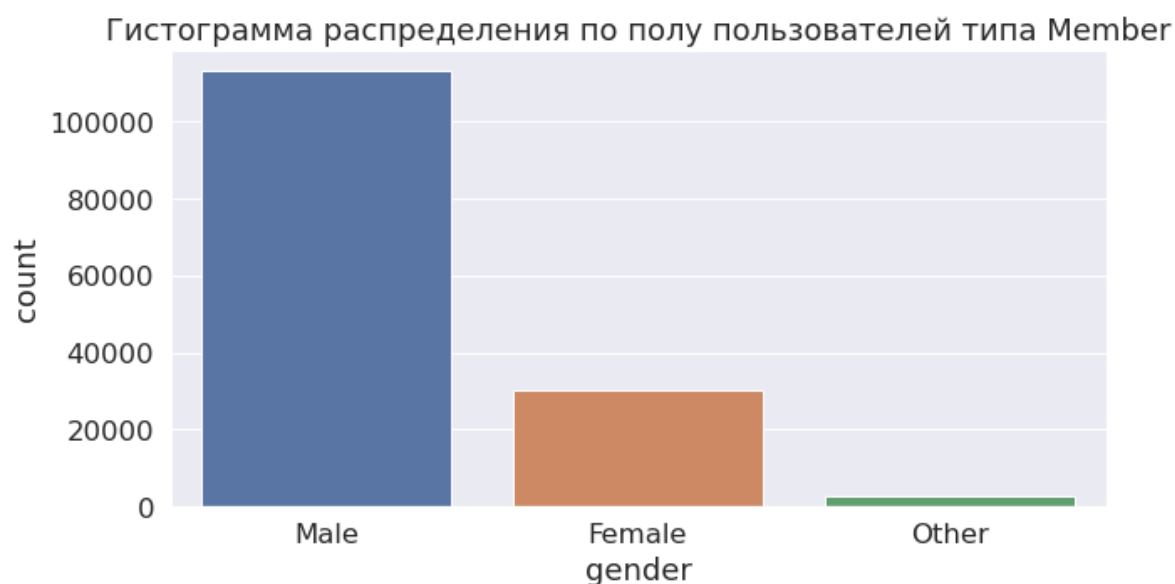
А вот Short-Term Pass Holder -пользователи велопрокатата на выходных пользуются уже чаще.

Посмотрим на распределение пользователей типа Member в зависимости от пола.

In [46]:

```
1 plt.figure(figsize=(10, 5))
2 plt.title('Гистограмма распределения по полу пользователей типа Member')
3 sns.countplot(x='gender', data=trips[trips['usertype'] == 'Member'])
4 plt.show()
```

executed in 988ms, finished 23:00:42 2020-04-26



При этом важно заметить, что пол пользователя для Short-Term Pass Holder подписки не определен, поэтому исследования на эту тему не возможны.

In [14]:

```
1 trips[trips.usertype != 'Member'].gender
```

executed in 91ms, finished 00:16:33 2020-05-21

Out[14]:

```
69      NaN
78      NaN
89      NaN
91      NaN
92      NaN
```

```
...
286852   NaN
286853   NaN
286854   NaN
286855   NaN
286856   NaN
```

Name: gender, Length: 89894, dtype: object

**Выводы:**

- Как видим, разные пользователи по-разному используют велопрокат. Так, пользователи типа Member катаются чаще и преимущественно небольшое время, при этом большая часть поездок приходится в начало (8 утра) и конец рабочего дня(17-18 вечера), а вот для Short-Term Pass Holder более характерны продолжительные поездки, о чем говорит и значение средних величин,

при чем совершаются они преимущественно днем, в период с 10 до 19 часов. Пользователя типа Member также очень редко катаются в выходные, чего нельзя сказать о Short-Term Pass Holder.