

```
In [48]: 1 import numpy as np
2 import pandas as pd
3 import scipy.stats as sps
4 import statsmodels as sm
5 from statsmodels.sandbox.stats.multicomp import multipletests
6
7 import matplotlib.pyplot as plt
8 import seaborn as sns
9 sns.set(font_scale=1.3)
10
11 %matplotlib inline
```

Данные о внебрачных отношениях

Данные используются для объяснения распределения времени между работой, временем, проведенным с супругом/супругой, и временем, проведенным с любовником/любовницей.

<http://www.statsmodels.org/stable/datasets/generated/fair.html>
(<http://www.statsmodels.org/stable/datasets/generated/fair.html>)

Исследуем, как каждый фактор влияет на долю времени, проведенного во внебрачных отношениях.

```
1 Number of observations: 6366
2 Number of variables: 9
3 Variable name definitions:
4
5     rate_marriage   : How rate marriage, 1 = very poor, 2 = poor, 3 = fair,
6                     4 = good, 5 = very good
7     age             : Age
8     yrs_married     : No. years married. Interval approximations. See
9                     original paper for detailed explanation.
10    children        : No. children
11    religious        : How religious, 1 = not, 2 = mildly, 3 = fairly,
12                     4 = strongly
13    educ             : Level of education, 9 = grade school, 12 = high
14                     school, 14 = some college, 16 = college graduate,
15                     17 = some graduate school, 20 = advanced degree
16    occupation       : 1 = student, 2 = farming, agriculture; semi-skilled,
17                     or unskilled worker; 3 = white-collor; 4 = teacher
18                     counselor social worker, nurse; artist, writers;
19                     technician, skilled worker, 5 = managerial,
20                     administrative, business, 6 = professional with
21                     advanced degree
22    occupation_husb  : Husband's occupation. Same as occupation.
23    affairs          : measure of time spent in extramarital affairs
24
25 See the original paper for more details.
```

Посмотрим на данные. Все переменные, кроме `affairs` являются категориальными, а переменная `affairs` --- вещественной.

```
In [3]: 1 data = pd.DataFrame(sm.datasets.fair.load().data)
2 data.head()
```

```
Out[3]:
```

	rate_marriage	age	yrs_married	children	religious	educ	occupation	occupation_husb	affairs
0	3.0	32.0	9.0	3.0	3.0	17.0	2.0	5.0	0.111111
1	3.0	27.0	13.0	3.0	1.0	14.0	3.0	4.0	3.230769
2	4.0	22.0	2.5	0.0	1.0	16.0	3.0	5.0	1.400000
3	4.0	37.0	16.5	4.0	3.0	16.0	5.0	5.0	0.727273
4	5.0	27.0	9.0	1.0	1.0	14.0	3.0	4.0	4.666666

В данных у 2/3 людей вообще не было внебрачных отношений

```
In [4]: 1 (data['affairs'] == 0).mean()
```

```
Out[4]: 0.6775054979579014
```

У таких людей все хорошо, и мы их не рассматриваем, так что просто удалим их.

```
In [5]: 1 data = data[data['affairs'] > 0]
        2 len(data)
```

```
Out[5]: 2053
```

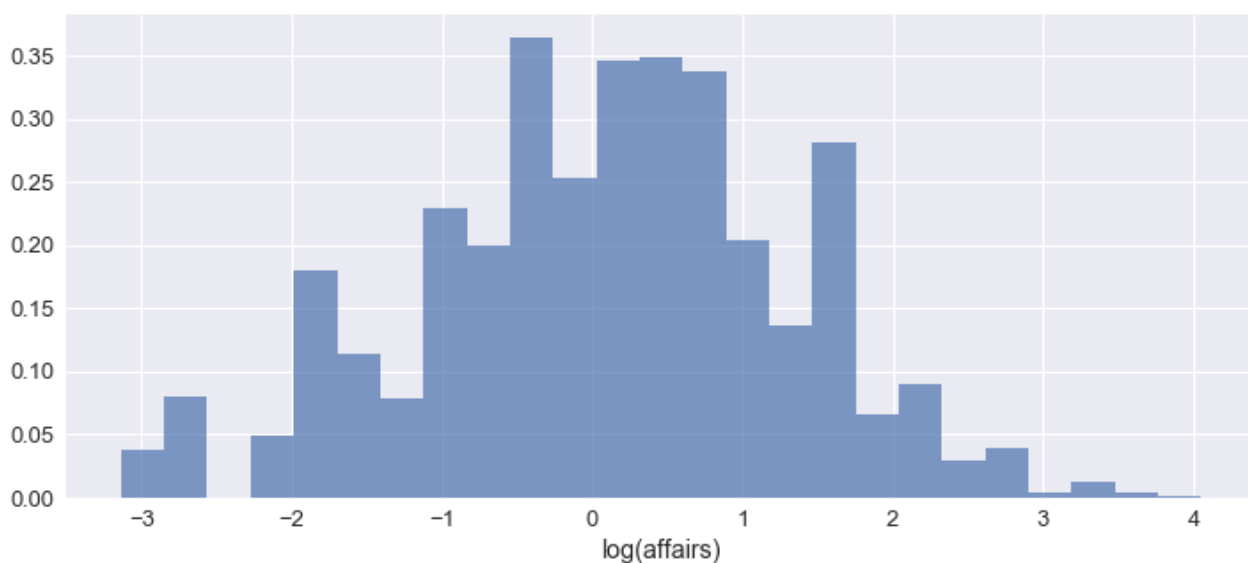
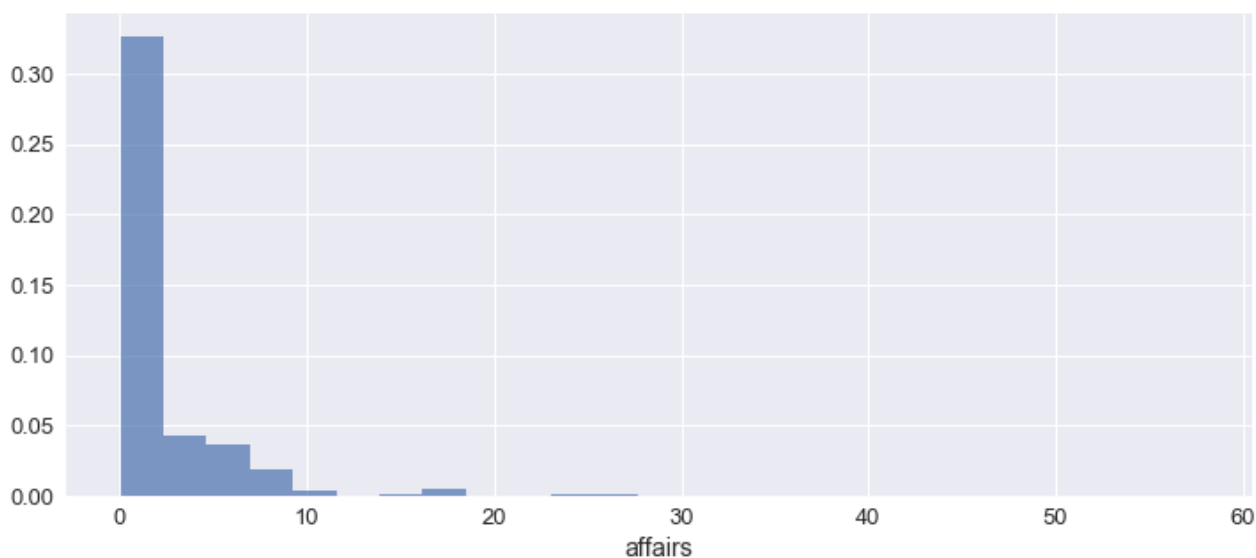
Описательные статистики по времени во внебрачных отношениях по всем людям

```
In [6]: 1 data['affairs'].describe()
```

```
Out[6]: count      2053.000000
mean         2.187243
std          3.437478
min          0.043478
25%          0.521739
50%          1.217391
75%          2.177776
max          57.599991
Name: affairs, dtype: float64
```

Посмотрим на гистограмму времени во внебрачных отношениях и на гистограмму логарифма этой величины

```
In [9]: 1 plt.figure(figsize=(12, 5))
2 plt.hist(data['affairs'], bins=25, alpha=0.7, normed=True)
3 plt.xlabel('affairs');
4
5 data['log(affairs)'] = np.log(data['affairs'])
6
7 plt.figure(figsize=(12, 5))
8 plt.hist(data['log(affairs)'], bins=25, alpha=0.7, normed=True)
9 plt.xlabel('log(affairs)');
```



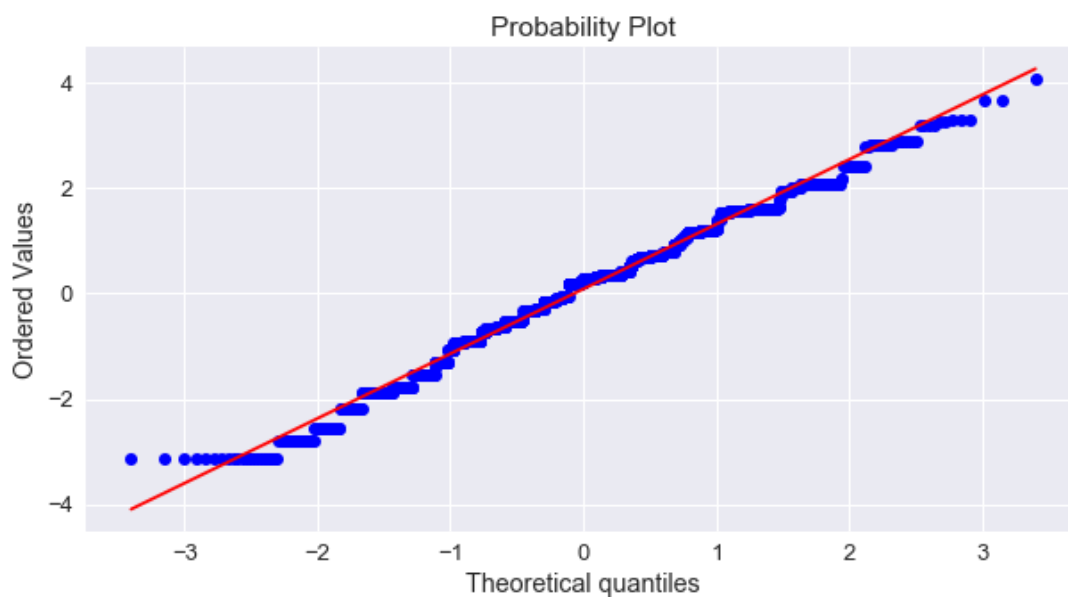
Логарифм времени папахивает нормальностью судя по гистограмме, но критерий Шапиро-Уилка отвергает ее

```
In [10]: 1 sps.shapiro(data['log(affairs)'])
```

```
Out[10]: (0.9899781942367554, 9.874644157914503e-11)
```

На QQ plot точки отдаленно расположены вдоль одной прямой

```
In [12]: 1 plt.figure(figsize=(10, 5))
2         ax = plt.subplot(111)
3         sps.probplot(data['log'affairs']), plot=ax);
```



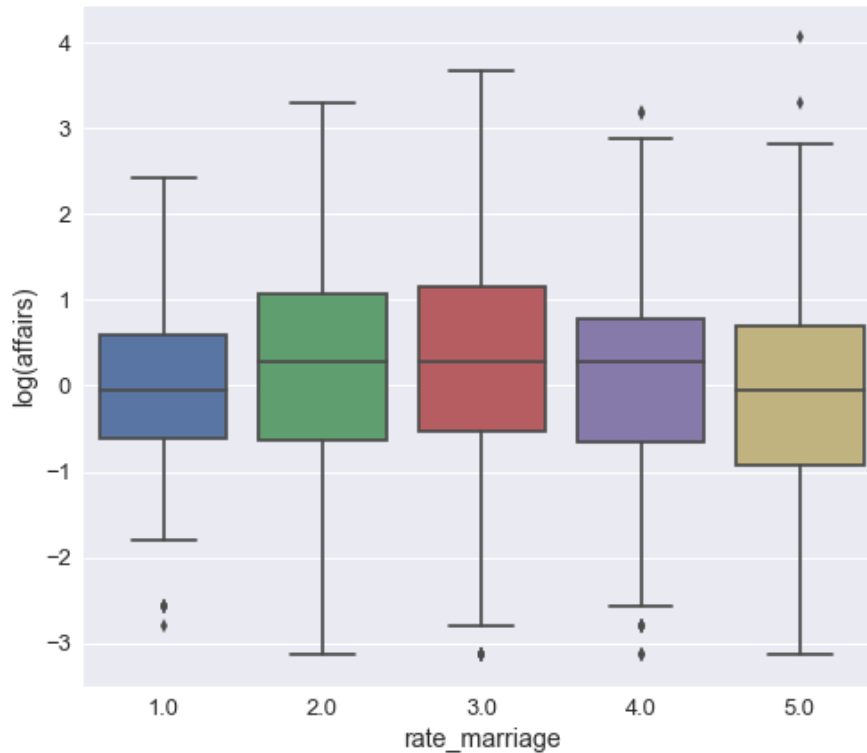
Анализ влияния факторов

```
In [45]: 1 def analyse_factor(factor_name):
2         print('Factor ' + factor_name)
3
4         gb = data['affairs'].groupby(by=data[factor_name])
5         samples = [np.array(group[1]) for group in gb]
6         kruskal_result = sps.kruskal(*samples)
7         print(kruskal_result)
8
9         plt.figure(figsize=(8, 7))
10        sns.boxplot(x=data[factor_name], y=data['log'affairs'])
11        plt.show()
12
13        return kruskal_result.pvalue
```

```
In [46]: 1 result = pd.DataFrame(columns=['factor', 'pvalue'])
2
3 for factor_name in ['rate_marriage', 'age', 'yrs_married',
4                     'children', 'religious', 'educ', 'occupation']:
5     pvalue = analyse_factor(factor_name)
6     result = result.append({'factor': factor_name, 'pvalue': pvalue}, ignore_index=True)
7     print('\n=====')
```

Factor rate_marriage

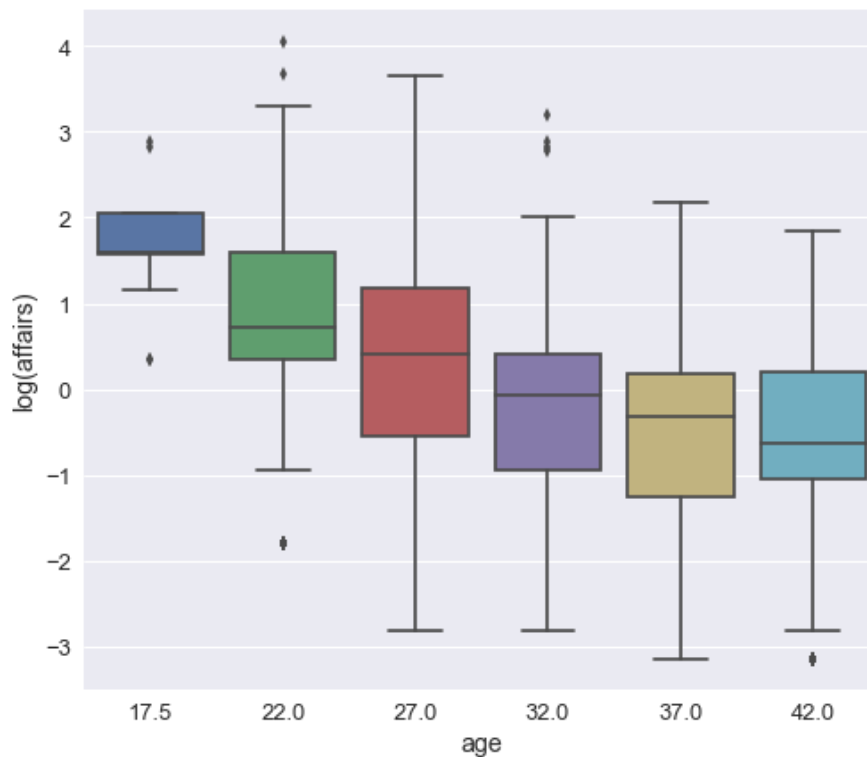
KruskalResult(statistic=18.79797401228851, pvalue=0.0008611183487432226)



=====

Factor age

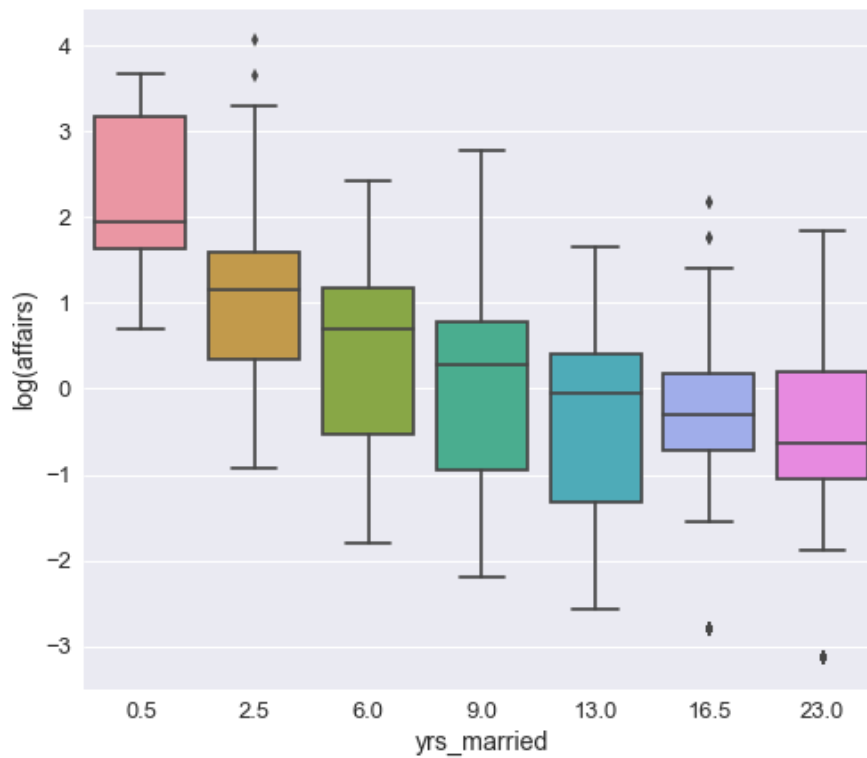
KruskalResult(statistic=366.92107249211216, pvalue=3.974891144400868e-77)



=====

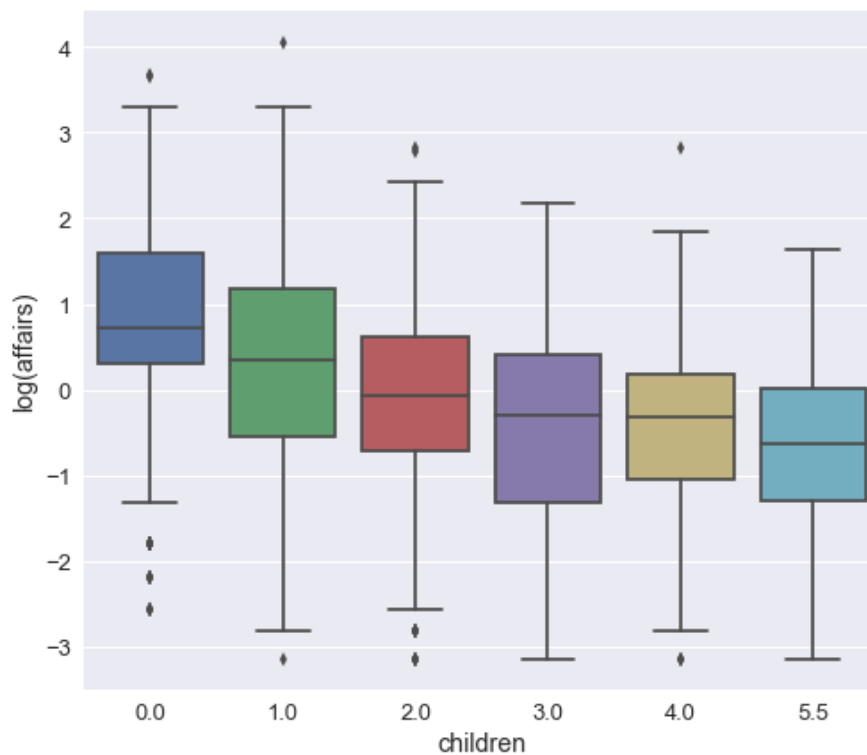
Factor yrs_married

KruskalResult(statistic=465.6566234684575, pvalue=2.0927681559863934e-97)



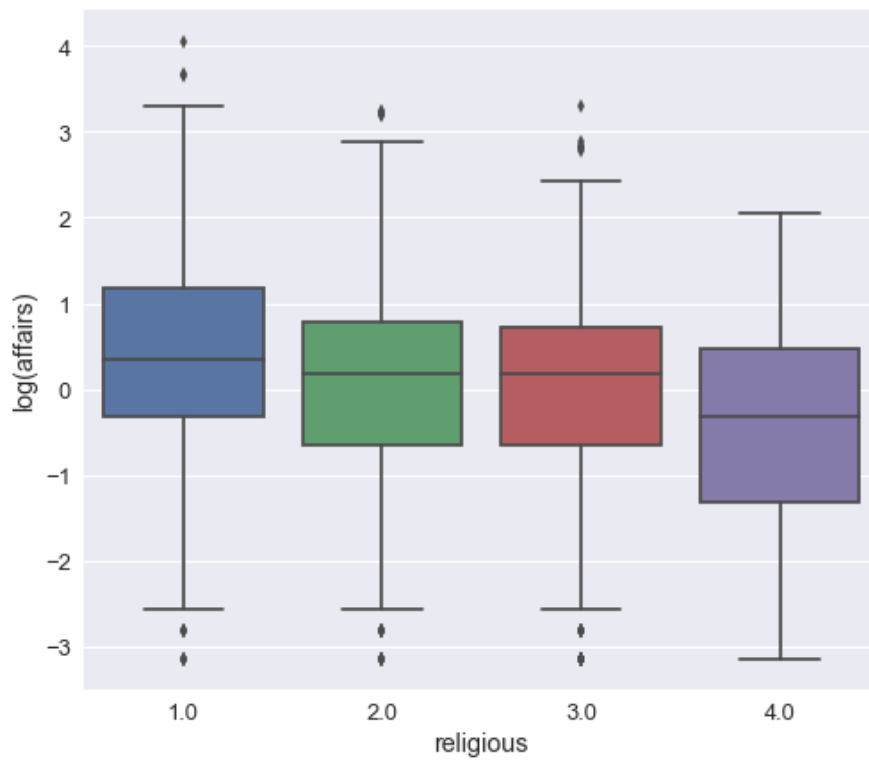
Factor children

KruskalResult(statistic=313.2218528260669, pvalue=1.437105588563059e-65)



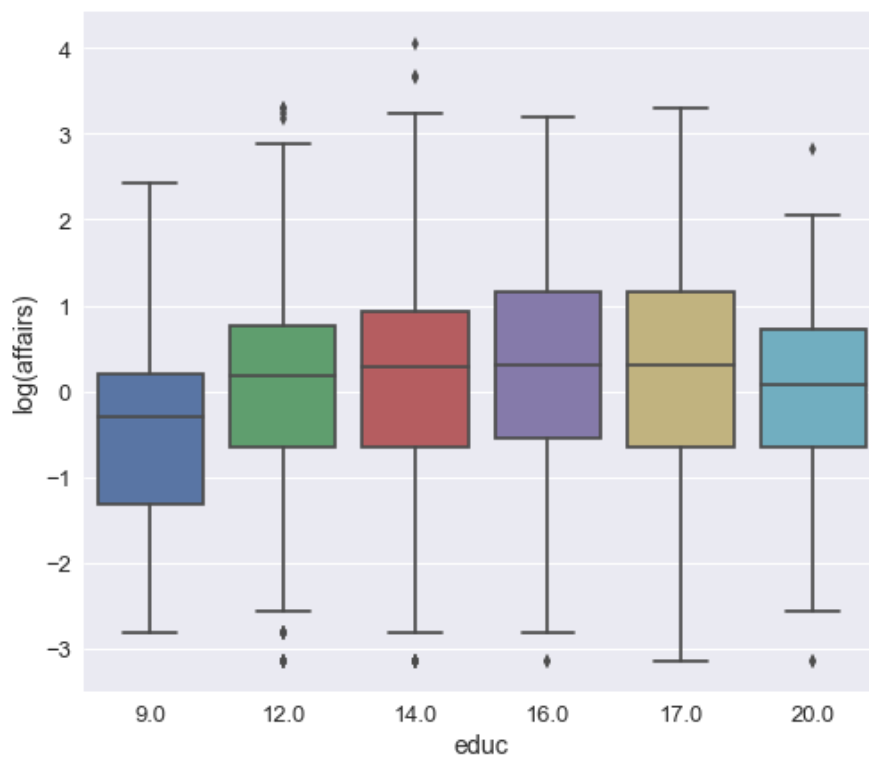
Factor religious

KruskalResult(statistic=44.566982906003965, pvalue=1.1436297209005018e-09)



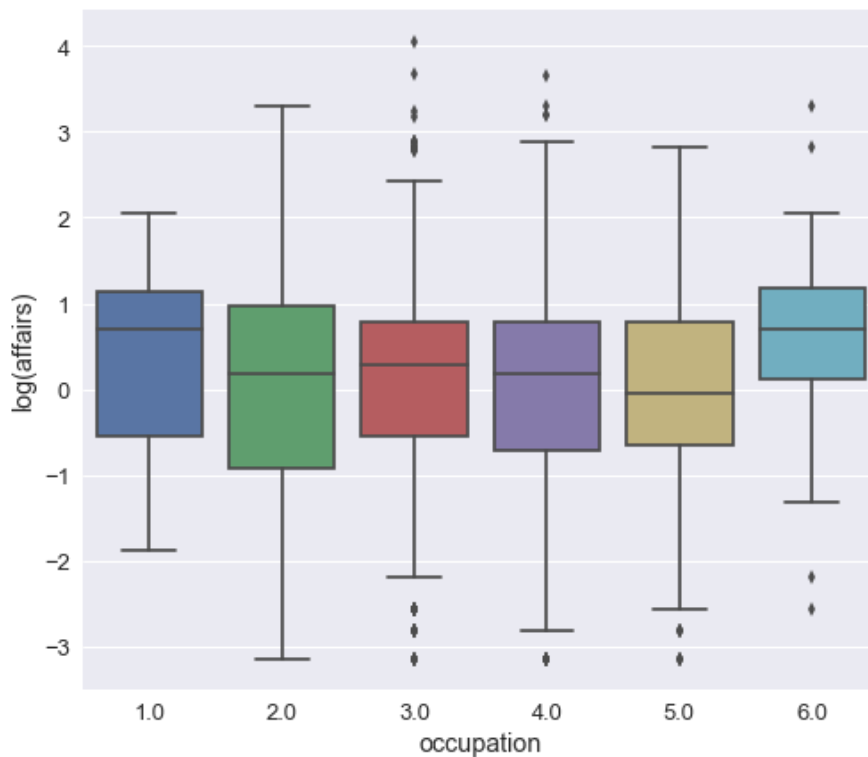
=====

Factor educ
 KruskalResult(statistic=12.061049042219755, pvalue=0.03396077556581535)



=====

Factor occupation
 KruskalResult(statistic=10.43824253852124, pvalue=0.06372848246987972)



=====

```
In [50]: 1 result['pvalue corrected'] = multipletests(result['pvalue'], method='holm')[1]
          2 result['reject'] = result['pvalue corrected'] < 0.05
          3
          4 result
```

```
Out[50]:
```

	factor	pvalue	pvalue corrected	reject
0	rate_marriage	8.611183e-04	2.583355e-03	True
1	age	3.974891e-77	2.384935e-76	True
2	yrs_married	2.092768e-97	1.464938e-96	True
3	children	1.437106e-65	7.185528e-65	True
4	religious	1.143630e-09	4.574519e-09	True
5	educ	3.396078e-02	6.792155e-02	False
6	occupation	6.372848e-02	6.792155e-02	False

Прикладная статистика и анализ данных, 2019

Никита Волков

<https://mipt-stats.gitlab.io/> (<https://mipt-stats.gitlab.io/>)