

# Анализ данных при помощи языка Python

Горин Н.А.

МГТУ им. Н.Э. Баумана

6 июля 2023 г.

## Аннотация

Для анализа данных были выбраны два исследования по наблюдению за морфологией Дарвиновых вьюрков (лат. *Geospizinae*), проведённые в 1975 г. и 2012 г. Питером и Розмари Грант (*Peter and Rosemary Grant*). Выбранная работа является классическим научным трудом, позволяющим проследить эволюцию живых организмов (вьюрков вида *Geospiza fortis* и *Geospiza scandens*) на коротком временном отрезке.

## Введение

В данной работе представлен статистический анализ данных, сделанный при помощи языка Python. При выполнении работы использовались библиотеки *matplotlib*, *numpy*, *pandas* и *seaborn*.

## Основная часть

Импортируем библиотеки и подключаем выборки:

```
1 import seaborn as sns
2 import pandas as pd
3 import numpy as np
4 import matplotlib.pyplot as plt
5
6 sns.set_theme()
7
8 beaks_1975 = pd.read_csv('1975.csv')
9 beaks_2012 = pd.read_csv('2012.csv')
```

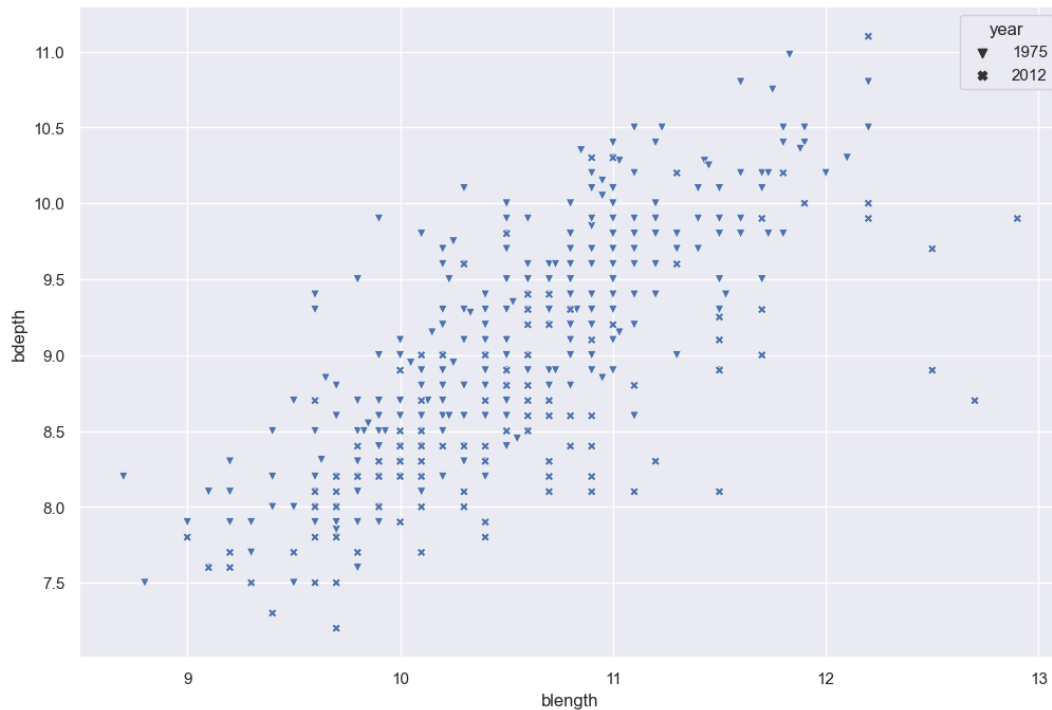
Редактируем файлы для более удобного анализа:

```
1 beaks_1975['year'] = "1975"
2 beaks_2012['year'] = "2012"
3 beaks_1975.rename(columns={'Beak depth, mm' : 'bdepth', 'Beak length, mm' :
    'blength'}, inplace=True)
4 data = pd.concat([beaks_1975, beaks_2012]).reset_index(drop=True)
5 fortis = data[data.species == 'fortis'].reset_index(drop=True)
6 scandens = data[data.species == 'scandens'].reset_index(drop=True)
```

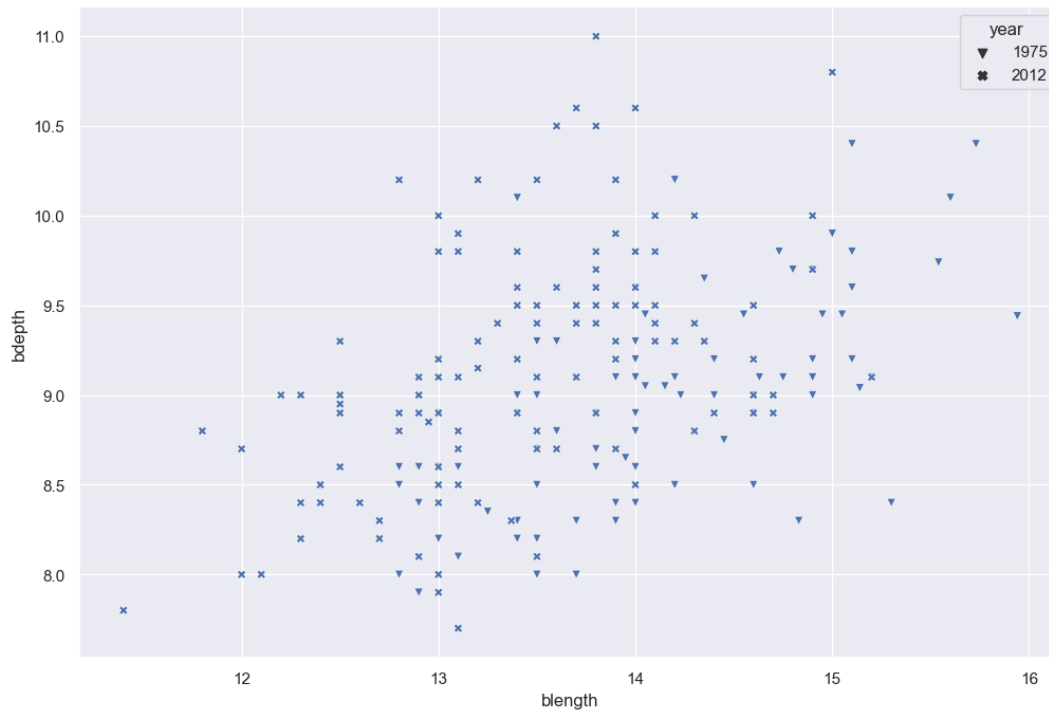
Маркируем данные и строим графики типа **ScatterPlot**:

```
1 markers = {'1975': "v", '2012': "X"}
2 fig = plt.figure(figsize=(12, 8))
3 sns.scatterplot(data=fortis, x='blength', y='bdepth', style='year', markers
    =markers)
4 plt.show()
5
6 markers = {'1975': "v", '2012': "X"}
7 fig = plt.figure(figsize=(12, 8))
8 sns.scatterplot(data=scandens, x='blength', y='bdepth', style='year',
    markers=markers)
9
10 plt.show()
```

Длина и ширина клюва (mm) вьюрков вида *fortis*



Длина и ширина клюва (mm) व्यूरков вида *scandens*

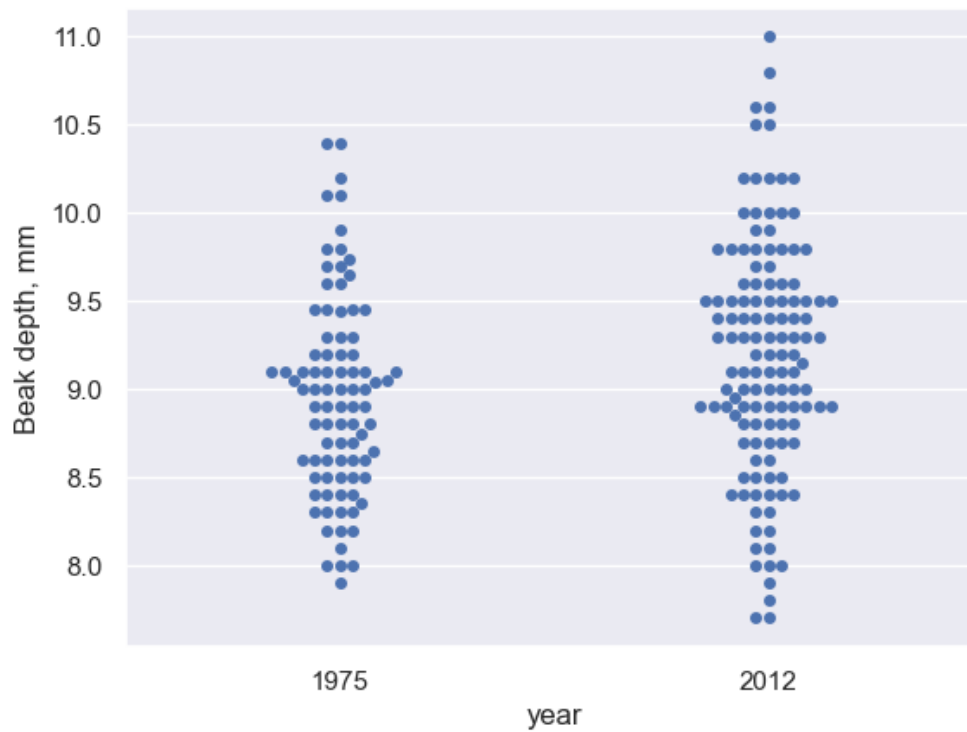


Из графиков видно заметное различие в параметрах клюва у व्यूरков вида *scandens*. Дальнейшие данные будем систематизировать по особям этого вида.

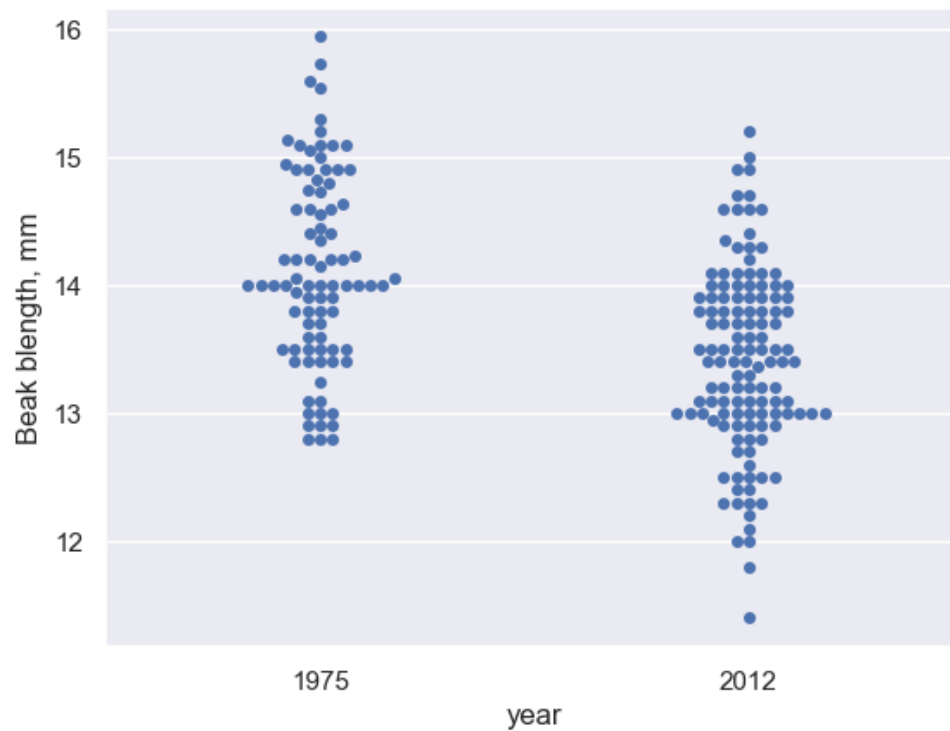
Разделим особей вида *scandens* по временным промежуткам проведённых исследований и построим графики типа **SwarmPlot**.

```
1 scandens_1975 = beaks_1975[beaks_1975['species'] == 'scandens']
2 scandens_2012 = beaks_2012[beaks_2012['species'] == 'scandens']
3 scandens_bdepth_1975 = scandens_1975['bdepth'].reset_index(drop=True)
4 scandens_bdepth_2012 = scandens_2012['bdepth'].reset_index(drop=True)
5
6 sns.swarmplot(scandens, x='year', y='bdepth')
7 plt.xlabel('year')
8 plt.ylabel('Beak depth, mm')
9
10 scandens_blength_1975 = scandens_1975['blength'].reset_index(drop=True)
11 scandens_blength_2012 = scandens_2012['blength'].reset_index(drop=True)
12 sns.swarmplot(scandens, x='year', y='blength')
13 plt.xlabel('year')
14 plt.ylabel('Beak blength, mm')
15
16 plt.show()
```

Ширина клюва выюров вида *scandens* за 1975 и 2012 гг.



Длина клюва выюров вида *scandens* за 1975 и 2012 гг.

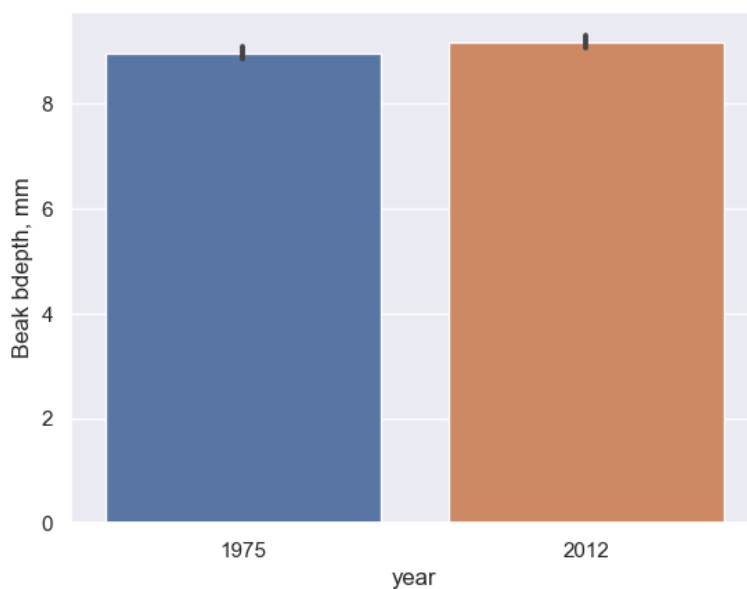


Построенные графики позволяют выявить явные статистические изменения в морфологии строения клюва вида *scandens* за 37 лет.

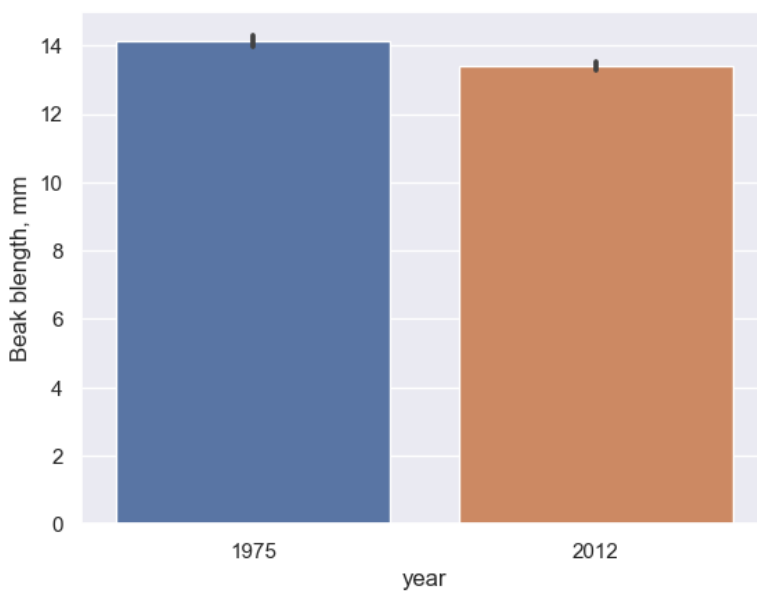
Выполним проверку результатов на графиках другого типа: **BarPlot**.

```
1 sns.barplot(scandens, x="year", y="bdepth")
2 plt.xlabel('year')
3 plt.ylabel('Beak bdepth, mm')
4
5 sns.barplot(scandens, x="year", y="blength")
6 plt.xlabel('year')
7 plt.ylabel('Beak blength, mm')
8
9 plt.show()
```

Ширина клюва выюрков вида *scandens* за 1975 и 2012 гг.



Длина клюва выюрков вида *scandens* за 1975 и 2012 гг.

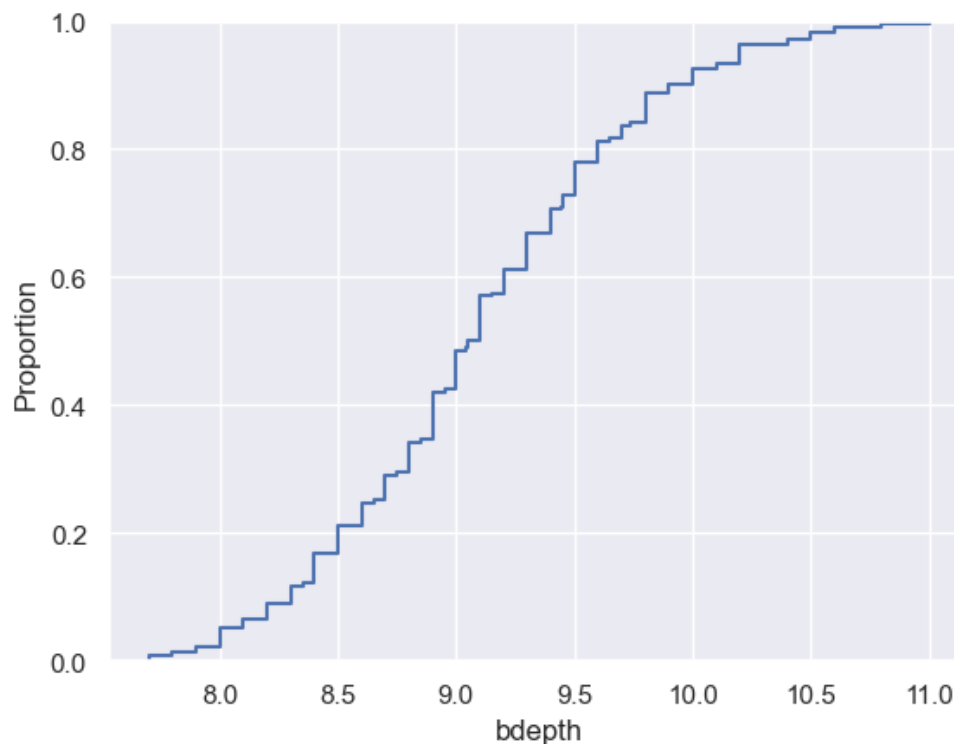


В завершение работы посмотрим, как менялась ширина и длина клюва вьюрков вида *scandens* в совокупности за 1975 и 2012 гг. Выполнить анализ помогут графики типа **ECDFPlot**.

Функция *ECDF* (эмпирическая кумулятивная функция распределения) представляет долю или количество наблюдений, попадающих ниже каждого уникального значения в наборе данных. Преимущество по сравнению с гистограммой или графиком заключается в том, что каждое наблюдение визуализируется напрямую, а это означает, что нет необходимости корректировать параметры группирования или сглаживания.

```
1 sns.ecdfplot(scandens, x="bdepth")
2 plt.xlabel('bdepth')
3 plt.ylabel('Proportion')
4
5 ns.ecdfplot(scandens, x="blength")
6 plt.xlabel('blength')
7 plt.ylabel('Proportion')
8
9 plt.show()
```

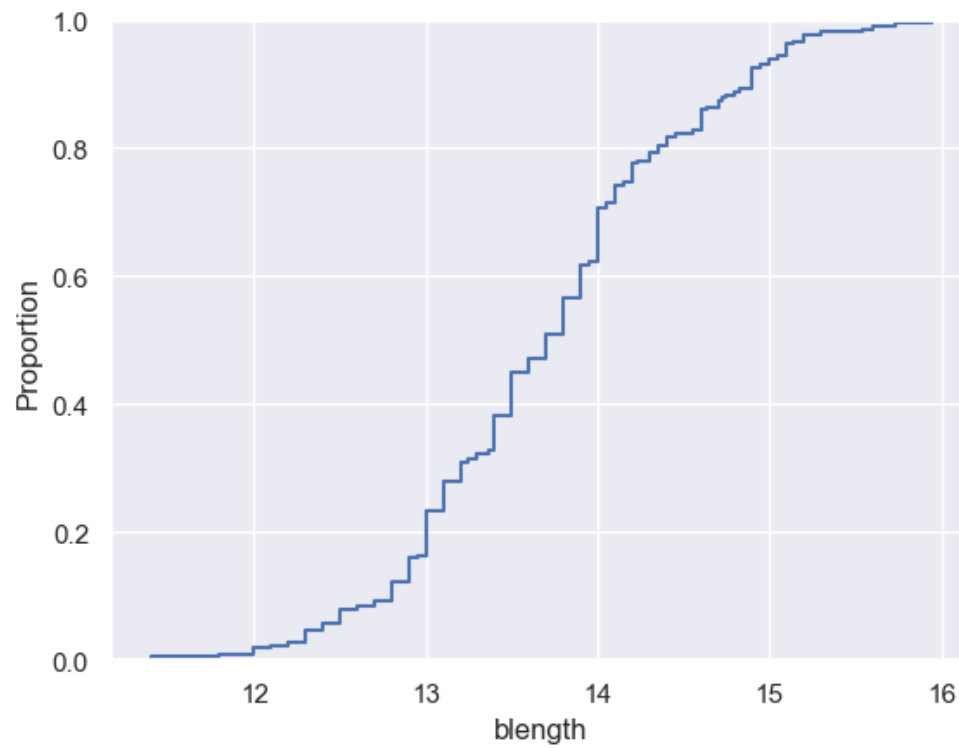
Ширина клюва вьюрков вида *scandens*



## Заключение

Проведённый статистический анализ является прямым свидетельством морфологической изменчивости вида *scandens*, произошедшими за

Длина клюва выюлков вида *scandens*



37 лет между исследованиями.

Результатом работы служит новое доказательство эволюционной теории Ч. Дарвина.