

# Разведочный Анализ Данных

Пилипенко А.О., Банников Д.А. гр.311, ВМК МГУ

## Постановка задачи

1. Необходимо провести Разведочный Анализ Данных (**Exploratory Data Analysis - EDA**) над имеющимися данными.
2. Написание программы, выполняющей Разведочный Анализ Данных.
  - a. Математическая интерпретация задачи и его решения.
  - b. Описание основных функций, используемых при программном решении.
  - c. Оформление кода в виде Python Notebook.
3. Результаты анализа и выводы оформляются в виде презентации.
  - a. Использование средств LaTeX и Beamer.
4. Описание кода программы, полученного в результате решения построенного примера.

## Описание задачи

После оглушительного успеха в освобождении Астапора, Миэрина и Юнка от власти работорговцев Дейенерис Бурерожденная открыла себе доступ к Летнему морю, а следовательно -- путь в Вестерос.

Для ведения войны с Семью Королевствами нужно оружие, а для оружия нужна сталь. Нет никаких сомнений в кузнечном искусстве Безупречных, однако поставщики стали не столь надежны.

Два основных поставщика стали - это Westeros Inc. и Harpy & Co. На протяжении нескольких месяцев мы покупаем сталь у обеих компаний, и каждая из них предлагает ощутимую скидку при заключении эксклюзивного договора на поставку.

Советник королевы Тирион Ланнистер знает о вашем умении принимать взвешенные рациональные решения и просит помощи в объективном решении вопроса о том, с какой из компаний следует заключить эксклюзивный договор на поставку стали.

У Тириона есть записи о производстве мечей каждым из кузнецов-безупречных, а также данные о количестве сломанных мечей в каждый из месяцев ведения боевых действий.

Вам дан CSV-файл с данными о производстве оружия и количестве единиц сломанного оружия за каждый месяц каждым из кузнецов.

Цель: ответ на вопрос: "С каким из поставщиков стали следует заключить договор?"

## Математическое решение задачи

Входные данные представляют собой набор строк формата:  
[ ДАТА\_ПРОИЗ, ДАТА\_ОТЧЕТА, ПРОИЗВЕДЕНО, ДЕФЕКТНЫЕ, ПОСТАВЩИК ]

где:

ДАТА\_ПРОИЗ - месяц, в котором была произведена продукция

(В нашем примере мы рассмотрим статистику на 6 месяцев производства)

ДАТА\_ОТЧЕТА - месяц, в котором были собраны данные о дефектности продукции, произведенной в месяце ДАТА\_ПРОИЗ

(Для участия данных о произведенных в последнем месяце элементов в анализе данных о качестве продукции компаний, в нашем примере отчеты будут производиться на один месяц дольше - 7 месяцев)

ПРОИЗВЕДЕНО - количество произведенной продукции в месяце ДАТА\_ПРОИЗ

ДЕФЕКТНЫЕ - количество обнаруженной дефектной продукции среди изделий, произведенных в месяце ДАТА\_ПРОИЗ, после проведения отчета в месяце ДАТА\_ОТЧЕТА

ПОСТАВЩИК - наименование компании-поставщика: Westeros Inc. или Harpy & C

Для проведения Разведочного Анализа Данных на наличие дефектности в производимой компаниями продукции, мы используем три экономических инструмента.

## Дефектность в среднем

Каждая компания характеризуется скоростью поломки ее продукции. Так как в явном виде мы не можем определить скорость становления продукции дефектной, мы по имеющимся данным можем построить статистику поломки продукции после каждого месяца эксплуатации.

Таким образом мы получим для всей продукции данные после одного месяца эксплуатации.

Для каждой строки из нашего набора данных:

Если ( ДАТА\_ОТЧЕТА - ДАТА\_ПРОИЗ = 1 ), то суммируем ДЕФЕКТНЫЕ в переменную суммирования СУМ\_1.

После двух месяцев эксплуатации:

Если ( ДАТА\_ОТЧЕТА - ДАТА\_ПРОИЗ = 2 ), то суммируем ДЕФЕКТНЫЕ в переменную суммирования СУМ\_2.

При этом надо понимать, что для произведенной на последнем месяце производства продукции статистика об эксплуатации после двух месяцев уже не будет учитываться.

Количество таких сумм определяется максимальным диапазоном эксплуатации:

МАКСИМУМ( ДАТА\_ОТЧЕТА - ДАТА\_ПРОИЗ ) = 6

Объединив данные по компаниям в единый график, мы получим базовое представление о том, как быстро продукция подвергается повреждению в сравниваемых нами компаниях, что станет для нас первым определителем надежности продукции каждой из интересующих нас корпораций.

## Совокупный срок службы

Одним из параметров продукции является ожидаемое время целостности продукции. При закупке продукции в больших объемах часто используют совокупный срок службы для получения более точных оценок реального объема получаемого товара.

Определить совокупный срок службы для каждой компании можно так: Предположим, что продукция была произведена в первом месяце и у нее не возникало дефектов все семь месяцев эксплуатации. Тогда продукция будет иметь совокупный срок службы:

$\text{ПРОИЗВЕДЕНО} * 7$ , где  $7 = \text{МАКСИМУМ}(\text{ДАТА\_ОТЧЕТА})$

Однако, часть данной продукции станет дефектной, предположим, в месяц  $\text{ДАТА\_ОТЧЕТА\_1}$ , потому вычитается величина:

$\text{ДЕФЕКТНЫЕ} * (7 - (\text{ДАТА\_ОТЧЕТА\_1} - 1))$

Т.е. если часть продукции сломается в первом же месяце, тогда же, когда ее произвели, то получится:

$(\text{ПРОИЗВЕДЕНО} - \text{ДЕФЕКТНЫЕ}) * 7$

Часть сломалась во втором месяце:

$\text{ПРОИЗВЕДЕНО} * 7 - \text{ДЕФЕКТНЫЕ} * 6$

Что можно интерпретировать так:

$(\text{ПРОИЗВЕДЕНО} - \text{ДЕФЕКТНЫЕ}) * 6 + \text{ПРОИЗВЕДЕНО} * 1$

Т.е. часть целой продукции прожила на один месяц дольше.

Учитывая, что производство происходит не только в первом месяце, а также повторив рассуждения выше, получим итоговую формулу:

Отдельно для каждого месяца  $\text{ДАТА\_ПРОИЗ}$  производства продукции

Суммируем по всем месяцам  $\text{ДАТА\_ОТЧЕТА}$  величины

$(\text{ПРОИЗВЕДЕНО} - \text{ДЕФЕКТНЫЕ}) * (8 - \text{ДАТА\_ОТЧЕТА})$

Эта величина средне Человеко-Часам в сфере трудоустройства описывает средний объем Продукто-Месяцев, предоставляемых компанией производителем. Совокупный срок службы является экономической интерпретацией объема закупленного товара. Таким образом мы получили второй инструмент сравнения интересующих нас фирм.

## Продукция без дефектов

Последним критерием для нас станут данные об объемах закупленной продукции, еще не ставшей дефектной, на каждый месяц эксплуатации:

Суммируем для каждого месяца  $\text{ДАТА\_ОТЧЕТА\_1}$  величины

$(\text{ПРОИЗВЕДЕНО} - \text{ДЕФЕКТНЫЕ})$

по всем месяцам  $\text{ДАТА\_ОТЧЕТА}$  с первого по  $\text{ДАТА\_ОТЧЕТА\_1}$

Полученные данные позволят нам определить количество имеющейся в результате всех закупок к моменту времени  $\text{ДАТА\_ОТЧЕТА\_1}$  целой продукции. При относительно малых скоростях порчи продукции данные величины показывают результаты всех закупок, произведенных с первого месяца, а в целом они определяют объем продукции, находящейся на складах покупателя, еще не подвергшейся поломке.

Таким образом, мы определили три статистических метода интерпретации данных, позволяющие получить аналитические сведения о надежности обеих компаний: Westeros Inc. и Harry & Co. Не является редкостью создание функций, по данным экономическим статистикам определяющих надежность компании в виде одного числа, например, при помощи весовых коэффициентов. Но мы предоставляем дальнейший анализ полученных данных пользователю.

Реализация анализа описываемого нами примера и полученные из него выводы оформлены в виде презентации.

## Основные функции

Разработка программы, реализующей полученные теоретические результаты и производящей графическое представление их пользователю, проходит в среде *Python Notebook*.

Основная задача программы - визуализация полученных нами результатов в виде двух типов графиков: стандартного с данными по месяцам и объемам, а так же Ящик с усами (boxplot).

Используемые библиотеки: *numpy*, *pandas*, *matplotlib.pyplot*, *seaborn*.

Библиотека *numpy* очень популярна, благодаря возможности выполнять различные математические операции с нестандартными для системы данными - массивами, списками и т.п. Функционал данной библиотеки интуитивно понятен, потому каждую функцию из данной библиотеки, используемую в нашем коде, мы не будем описывать сейчас.

Библиотека *pandas* используется для анализа данных, в частности, она позволяет нам работать с файлами CSV. Для этого используются функции:

**read\_csv( adr )** - Чтение данных из CSV файла, расположенного по адресу *adr*.

**DataFrame( )** - Конструктор объекта DataFrame (основной тип данных в *pandas*)

**X.groupby( Y )** - Объединение всех строк с одинаковыми элементами в столбцах Y в одну.

**X.mean( )** - Округление всех данных в объекте.

**X.transpose( )** - Транспонирование (разворот) данных (столбцы < - > строки).

**X.columns** - Доступ к списку названий столбцов

*Matplotlib.pyplot* - упрощенный интерфейс библиотеки *matplotlib* создания интерактивных графиков программными средствами. Используемые нами функции:

**title( X )** - Создание заголовка графика - X.

**show( )** - Отображение поля графиков.

**X.plot( )** - Построение стандартного графика x на y по данным X.

**Ax.set\_yticklabels( Y )** - выставляет на оси y (набора данных об осях - Ax) точечные пометки по списку Y.

**Ax.set\_xlim( X )** - Выставляет границы оси x (набора данных об осях - Ax) по промежутку X.

Библиотека *seaborn* является расширением библиотеки *matplotlib* для визуализации статистических данных. Нас интересует только одна функция:

**boxplot( Data )** - Построение графика Ящик с усами (boxplot) по данным Data.

## Код программы

Далее мы предоставляем текст программы, написанной нами для решения описанной выше задачи. Среда разработки - *Python Notebook*.

\*\*\*

```
1. import pandas as pd
2. import matplotlib.pyplot as plt
3. import numpy as np
4. import seaborn as sns
5.
6. pdata = pd.read_csv('production-data.csv', delimiter=',')
7. pdata['new'] = pdata['report.date'] - pdata['production.date']
8. gb = pdata.groupby(['supplier', 'new'])['defects']
9. gb = gb.mean()
10. harp_co = np.array(gb['harpy.co'][1:])
11. west_inc = np.array(gb['westeros.inc'][1:])
12.
13. d = (np.array([np.array(gb['harpy.co'][1:]),
14.                  np.array(gb['westeros.inc'][1:]))]).transpose()
14. df = pd.DataFrame(data=d)
15. df.columns = ['harpy.co', 'westeros.inc']
16. ax = df.plot(kind='barh', stacked=False)
17.
18. ax.set_yticklabels(['1 month', '2 month', '3 month', '4 month', '5 month', '6 month'])
19. ax.set_xlim([0, 20])#пределы для оси x
20. plt.title('Average number of defects per month')
21. plt.show()
22.
23. #После строиться график box
24. d = {'harpy.co': harp_co, 'westeros.inc': west_inc}
25. df = pd.DataFrame(data=
26. plt.title('Average number of defects per month')
27. sns.boxplot(data=df)
28. plt.show()
29.
30. sum_def_harp = np.zeros(6)
31. sum_def_west = np.zeros(6)
32. test_harp = np.zeros(6)
33. test_west = np.zeros(6)
34.
35. pdata = pd.read_csv('production-data.csv', delimiter=',', header=None)
36. pdata = np.array(pdata[1:])
37. n, m = np.shape(pdata)
38.
39. for i in range(n):
```

```

40. t_1 = float(pdata[i, 3]) - float(pdata[i, 4]) #produce - deffects
41. t_2 = t_1 * (8 - float(pdata[i, 2]))#float(pdata[i, 3]) * (8 - float(pdata[i, 2]))
42. #- float(pdata[i, 4]) * (8 - float(pdata[i, 2])) #how longe was working blade
43. if pdata[i, 5] == 'harpy.co':#fill if harpy.co
44.     j = int(pdata[i,1]) - 1
45.     test_harp[j] += t_1
46.     sum_def_harp[j] += t_2
47. if pdata[i, 5] == 'westeros.inc':#fill if westeros.inc
48.     j = int(pdata[i,1]) - 1
49.     test_west[j] += t_1
50.     sum_def_west[j] += t_2
51.
52. d = (np.array([sum_def_harp, sum_def_west])).transpose()
53. df = pd.DataFrame(data=d)
54. df.columns = ['harpy.co', 'westeros.inc']
55.
56. ax = df.plot(kind='barh', stacked=False)
57. ax.set_yticklabels(['1 month', '2 month', '3 month', '4 month', '5 month', '6 month'])
58. ax.set_xlim([0, 40000])#пределы для оси x
59. plt.title('Aggregate service life')
60. plt.show()
61.
62. #После строится график box
63. plt.title('Aggregate service life')
64. d = {'harpy.co': sum_def_harp, 'westeros.inc': sum_def_west}
65. df = pd.DataFrame(data=d)
66. sns.boxplot(data=df)
67. plt.show()
68.
69. d = (np.array([test_harp, test_west])).transpose()
70. df = pd.DataFrame(data=d)
71. df.columns = ['harpy.co', 'westeros.inc']
72. ax=df.plot(kind='barh', stacked=False)
73. ax.set_yticklabels(['1 month', '2 month', '3 month', '4 month', '5 month', '6 month'])
74. ax.set_xlim([0,6000])#пределы для оси x
75. plt.title('Summary number of swords without defects')
76. plt.show()
77.
78. #После строиться график box
79. plt.title('Summary number of swords without defects')
80. d = {'harpy.co': test_harp, 'westeros.inc': test_west}
81. df = pd.DataFrame(data=d)
82. sns.boxplot(data=df)
83. plt.show()

```

\*\*\*