

Московский авиационный институт  
(национальный исследовательский университет)

Факультет информационных технологий и прикладной  
математики

Кафедра вычислительной математики и программирования

Лабораторная работа №0 по курсу «Искусственный интеллект»

Студент: В. И. Лобов  
Преподаватель: С. Х. Ахмед  
Группа: М8О-306Б  
Дата:  
Оценка:  
Подпись:

Москва, 2019

# Лабораторная работа №0

**Задача:** Требуется сформировать/получить два набора данных соответствующие следующим критериям:

1. Один из датасетов должен представлять собой корпус документов. Язык, источник и тематика произвольна
2. Второй датасет должен содержать категориальные, количественные признаки. Для данного датасета определить предсказываемые признаки (для задачи регрессии и классификации). Если такого признака нет, спроектировать

Данные датасеты будут в дальнейшем использованы в оставшихся лабораторных работах.

По каждому датасету построить распределения признаков (в случае корпуса документов – построить распределение слов) и объяснить имеющуюся картину. Вычислить статистические характеристики признаков. Обнаружить и решить возможные проблемы с данными. Если решить данную проблему невозможно, объяснить почему

**Требования:**

1. Датасеты должны быть уникальны
2. Исходный код должен быть написан в одном код стайле
3. Должен быть указан источник данных

**Выбранные датасеты:**

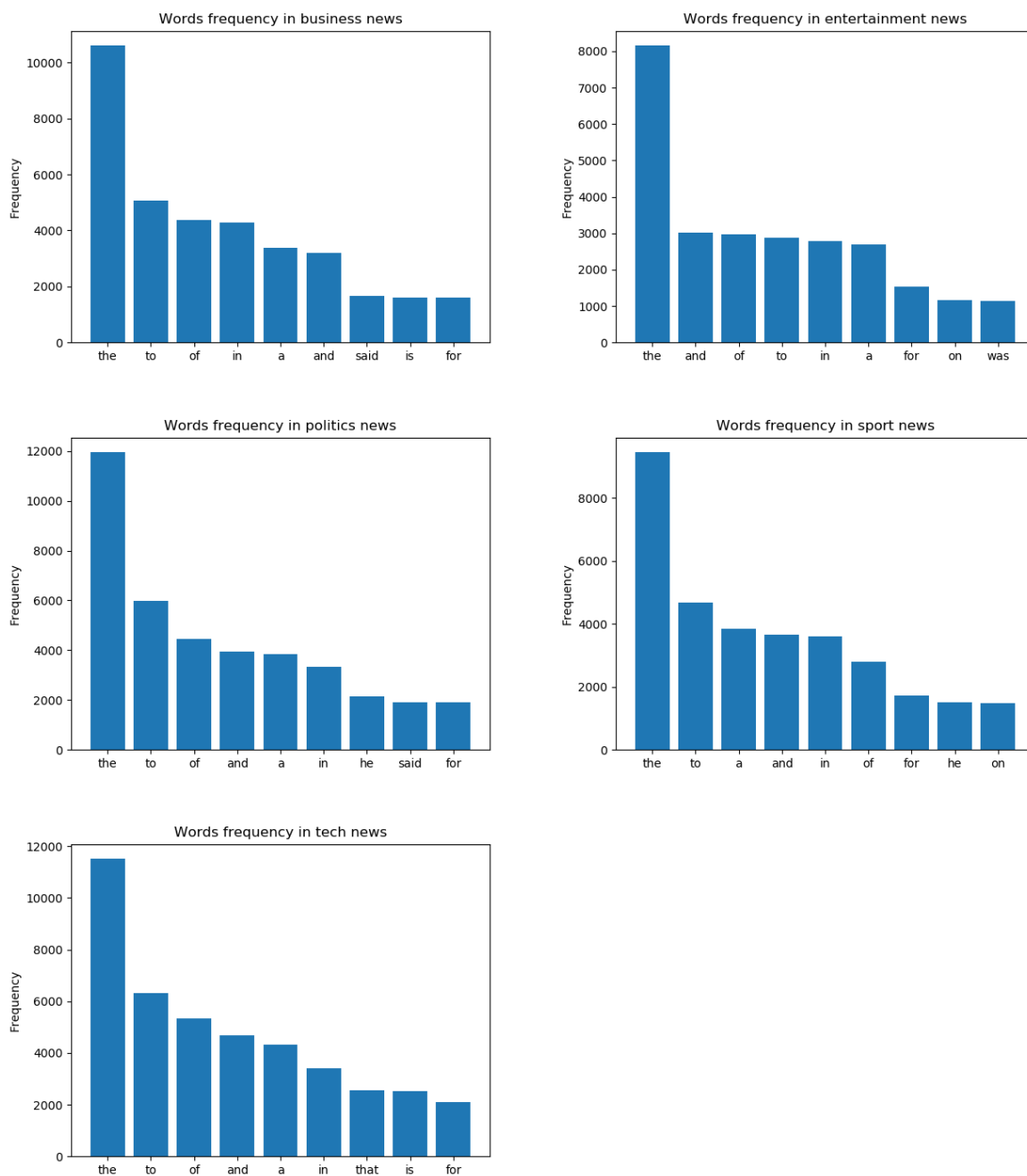
1. Корпус документов - статьи BBC по 5 жанрам: business, entertainment, politics, sport, tech  
<http://mlg.ucd.ie/datasets/bbc.html>
2. Данные котировок акций американских компаний(NYSE, NASDAQ)  
<https://www.kaggle.com/borismarjanovic/price-volume-data-for-all-us-stocks-etfs>

## 1 Описание

**Корпус документов.**

Так как стоит задача построения распределения слов в исходных текстах, необходимо устранить такие проблемы данных, как наличие знаков препинания и заглавных символов - выберем все вхождения знаков препинания в текст и приведём его к

нижнему регистру. Построим распределение слов в полученном обработанном тексте для каждого жанра, а затем выберем 10 наиболее встречающихся слов - получим следующую картину:



Проанализировав полученные результаты, можно сделать вывод, что вне зависимости от жанра текста наиболее популярными словами будут короткие (3-4 буквы)

слова-связки, а также слова, обозначающие косвенную/прямую речь (said), что свойственно для текстов типа новостей.

Листинг 1: src/bbc.py

```
1 import matplotlib.pyplot as plt
2 from os import listdir
3 import numpy as np
4
5 curDir = "C:\\Users\\pkmixer\\Downloads\\AI\\bbc\\"
6
7 for directory in listdir(curDir):
8     dictionary = {}
9     frequency = []
10    for file in listdir(curDir + directory):
11        with open(curDir + directory + "\\ " + file, "r") as f:
12            for line in f.readlines():
13                line = line.lower()
14
15                chars = [".", ",", "\'"]
16                for char in chars:
17                    line = line.replace(char, "")
18
19                line = line.split()
20
21                for word in line:
22                    if word not in dictionary:
23                        dictionary[word] = 0
24                    else:
25                        dictionary[word] += 1
26
27    for word in dictionary:
28        frequency.append((word, dictionary[word]))
29
30    frequency.sort(key=lambda x: x[1])
31
32    X = []
33    Y = []
34    for word in frequency[:-10:-1]:
35        X.append(word[0])
36        Y.append(word[1])
37
38    y_pos = np.arange(len(X))
39    plt.bar(y_pos, Y)
40    plt.xticks(y_pos, X)
41    plt.ylabel("Frequency")
42    plt.title("Words frequency in {} news".format(directory))
43    plt.savefig(curDir + "..\\bbc_pics\\" + directory + ".png")
44    plt.close()
```

### Данные котировок акций.

Исходные данные о котировках акций за некоторый промежуток времени представлены в виде отдельных файлов - каждый файл характеризует конкретную акцию. Можно соединить все акции в один .csv файл, при этом общий размер увеличится. Если сделать предположение, что котировки одной акции не зависят от котировок других (а так в большинстве случаев и есть), то можно рассматривать каждый файл, как отдельный датасет.

При рассмотрении датасета выяснилось, что некоторые файлы пусты, эти акции необходимо удалить, так как они не несут никакой смысловой нагрузки.

Листинг 2: src/a.py

```
1 | import os
2 | curDir = "C:\\Users\\pkmixer\\Downloads\\AI\\Stocks\\"
3 |
4 | for fileName in os.listdir(curDir):
5 |     if os.stat(curDir + fileName).st_size is 0:
6 |         print(fileName)
7 |         os.remove(curDir + fileName)
```

Каждая акция содержит 7 признаков:

1. Date - день котировки
2. Open - цена на открытии биржи
3. High - наибольшая цена за день
4. Low - наименьшая цена за день
5. Close - цена на закрытии биржи
6. Volume - суммарное количество операций за день
7. OpenInt(Open Interest) - сумма открытых позиций

В случае акций, величина Open Interest всегда равна 0 (имеет ненулевое значение для фьючерсов и других срочных контрактов), поэтому её стоит исключить из обработки данных. Volume - не слишком информативный признак, так как на предсказание цены акции влияет её история и экономическая ситуация, но никак не количество торгуемых акций, поэтому её так же можно исключить.

Построим распределение признаков акций, найдём среднее, дисперсию, медиану и некоторые квантили и изобразим наиболее "интересные" из них.

Листинг 3: src/statistics.py

```

1 import matplotlib.pyplot as plt
2 import pandas as pd
3 from pandas.plotting import scatter_matrix
4 import os
5
6
7 curDir = "C:\\Users\\pkmixer\\Downloads\\AI\\Stocks\\"
8
9 names = {}
10 indicesDir = "C:\\Users\\pkmixer\\Downloads\\indices\\"
11
12 for fileName in os.listdir(indicesDir):
13     with open(indicesDir + fileName, "r") as indiceName:
14         indiceName.readline()
15         for line in indiceName:
16             pair = line.strip().split("\t")
17             if len(pair) == 2:
18                 ticker, name = pair
19                 names[ticker.lower()] = name
20
21 skipColumns = ["Volume", "OpenInt"]
22
23 for fileName in os.listdir(curDir):
24     data = pd.read_csv(curDir + fileName, usecols=lambda x: x not in skipColumns)
25
26     index = fileName.split(".")[0]
27     if not names.get(index):
28         names[index] = index
29
30     scatter_matrix(data, alpha=0.2)
31     plt.suptitle(names[index])
32     plt.savefig(curDir + "pics\\" + fileName + ".png")
33     plt.close()

```

Листинг 4: src/stats.py

```

1 import matplotlib.pyplot as plt
2 import pandas as pd
3 import os
4
5
6 curDir = "C:\\Users\\pkmixer\\Downloads\\AI\\Stocks\\"
7
8 names = {}
9 indicesDir = "C:\\Users\\pkmixer\\Downloads\\indices\\"
10
11 for fileName in os.listdir(indicesDir):

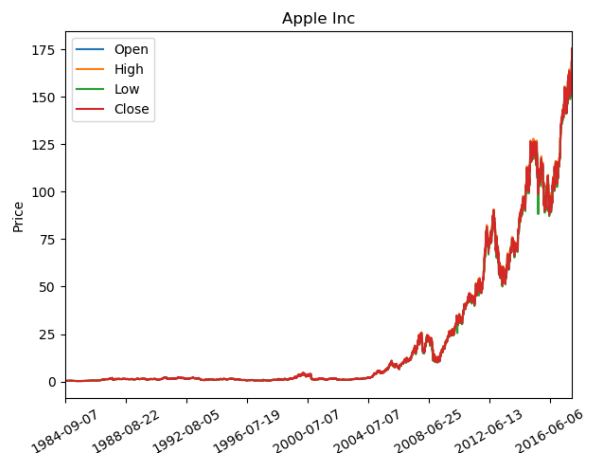
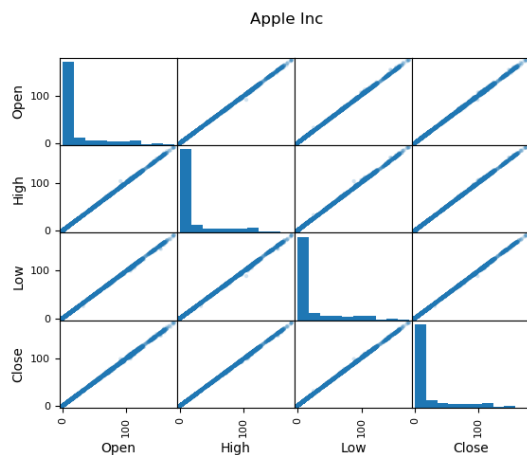
```

```

12     with open(indicesDir + fileName, "r") as indiceName:
13         indiceName.readline()
14         for line in indiceName:
15             pair = line.strip().split("\t")
16             if len(pair) == 2:
17                 ticker, name = pair
18                 names[ticker.lower()] = name
19
20 skipColumns = ["Date", "OpenInt", "Volume"]
21
22 for fileName in os.listdir(curDir):
23     data = pd.read_csv(curDir + fileName, usecols=lambda x: x not in skipColumns)
24
25     index = fileName.split(".")[0]
26     if not names.get(index):
27         names[index] = index
28
29     with open(curDir + "..\\\" + fileName + ".stat", "w") as file:
30         file.write(data.describe().to_csv())

```

Классический пример - акции Apple.



Метод `scatter_matrix` строит распределение признаков, а также зависимость одного признака от другого. Как видно из графика котировок, акция в среднем росла, что и отражено в виде прямой с углом наклона  $\frac{\pi}{4}$ .

Числовые характеристики:

Листинг 5: `src/aapl.csv`

```

1 Name,Open,High,Low,Close
2 count,8364.0,8364.0,8364.0,8364.0
3 mean,22.284350151841224,22.49586662362506,22.054243920373025,22.281018027259687
4 std,37.76340245539372,38.05773267058666,37.44743243649739,37.76446946661696
5 min,0.23305,0.23564000000000002,0.23051,0.23051

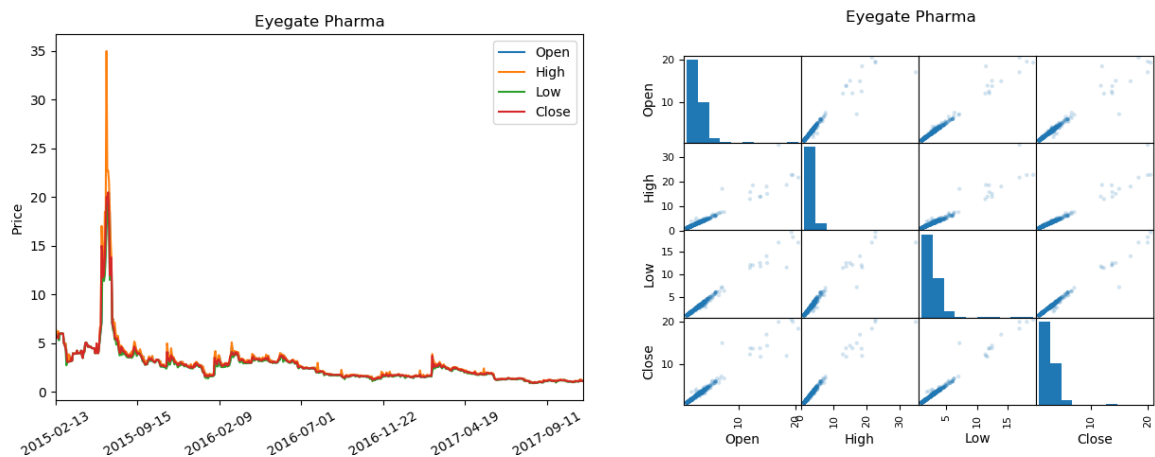
```

```

6 | 25%,1.1371,1.1642,1.1128,1.1371
7 | 50%,1.6328,1.6634,1.6006,1.62825
8 | 75%,23.739,23.930500000000002,23.335749999999997,23.6945
9 | max,175.11,175.61,174.27,175.61

```

Более интересная картина наблюдается с акциями Eyegate Pharma:



Характерный пик в стоимости акций отражается на зависимости одного признака от другого, а также в статистических характеристиках.

Листинг 6: src/eyeg.csv

```

1 | Name,Open,High,Low,Close
2 | count,645.0,645.0,645.0,645.0
3 | mean,2.78324992248062,2.9571705426356587,2.6390708527131785,2.773861395348837
4 | std,2.1770594051018417,2.7123703578832044,1.9780161518489183,2.2149790485909358
5 | min,0.919,0.95,0.9,0.93
6 | 25%,1.63,1.69,1.57,1.62
7 | 50%,2.43,2.49,2.36,2.42
8 | 75%,3.28,3.41,3.15,3.2648
9 | max,20.45,35.0,19.25,20.5

```

В данном датасете можно попробовать, например, предсказать цену акции на открытии торгов по данным, полученным за другой период, или максимальное падение акции за день в течение торгов.