# Statistics Guided Exercise

For using statistics with Python, we will be looking at the Pandas library. Pandas itself is built on top of another library, NumPy, and both have their own data structures. In this exercise, we will go over these data structures, and introduce you to Bokeh, which is a visualisation library you will be using in this exercise and the next for graphs and charts.

## Pandas

Pandas is a Python library for data analysis in Python. It providies some useful functions and data structures for the collection and analysis of data. In particular, we will be making use of the **DataFrame** and **Series** classes.

A `DataFrame` object represents data in a series of rows (individual observations of data) and columns (features or variables) within those data. Each of those rows and columns can be extracted, and they then become a `Series`. We will work through an example to illustrate these concepts.

### Importing Data

There are convenience functions to import data, such as **read_json** and **read_csv** which, as their names suggest, will import data which is already in a particular format. For this example, we will import data from the MongoDB database we used in the exercise last week.

For this example, we will import the first 1000 documents (*MongoDB stores data records as BSON documents. BSON is a binary representation of JSON documents, though it contains more data types than JSON*), in the UK collection into a Pandas `DataFrame`. Run the cell below

In [1]:

```python
# Convention is to import numpy and pandas with abbreviated names
# This means that instead of using pandas.read_csv, you would use pd.read_csv
import numpy as np
import pandas as pd #넘파이의 상위 라이브러리

from bokeh.io import output_notebook, show #파이썬 시각화 응용프로그램 보케
from bokeh.charts import * #차트를 하위 라이브러리로 만들고 보케의 구성요소가 됩니당

# Import PyMongo, so that we can query some data
# 'mongodb://cpduser:M13pV5woDW@mongodb/health_data' is the location of the data we

from pymongo import MongoClient #몽고DB
client = MongoClient('mongodb://cpduser:M13pV5woDW@mongodb/health_data')
db = client.health_data

cursor = db.uk.find({}).limit(1000) #할당을 할 때 결과에 제한(1000개 넘지 말것) 을 걸었음
                          #빈 딕셔너리에
# Unfortunately, Pandas does not support PyMongo objects for import, so we need to (
listy = list(cursor) #변수 커서의 유형을 리스트로 설정, 판다스를 써서 쓸 수 있게 하고 있음

# Create a Pandas DataFrame with the list object as a parameter
first_1000 = pd.DataFrame(listy) #리스트인 리스티를 데이터 프레임 형식으로 바꿈
print(first_1000)
```

|   | AddressLine1 | AddressLine2 | Address Line3 |
|---|---|---|---|
| 0 | NaN | 16a Adelphi Street | Pr eston |
| 1 | NaN | 24 Fylde Road | Pr eston |
| 2 | NaN | 119a Church Street | Pr eston |
| 3 | NaN | 30-34 Holmrook Road | Pr eston |
| 4 | NaN | NaN | NaN |
| 5 | NaN | 109-113 Avenham Lane | Pr eston |
| 6 | NaN | NaN | NaN |
| 7 | NaN | 515 Blackpool Road | Pr eston |
| 8 | NaN | 27 Meadow Street | Pr eston |

Now we have our imported data in a **DataFrame object**. Like any other Python object, it has a collection of attributes and methods which we can use. We will go over some here, but see the documentation (http://pandas.pydata.org/pandas-docs/stable/generated/pandas.DataFrame.html) for a full list. We'll start by seeing what the data looks like by calling the `head()` function on the data:

In [29]:

```
# Filtering the data in the DataFrame to only return rows where RatingValue < 3
first_1000[first_1000['RatingValue'] < 3]
#결과 값이 레이팅벨류가 0, 1, 2 인 것만 출력
```

| | | | | | | |
|---|---|---|---|---|---|---|
| **292** | NaN | 49 Blackpool Road | Preston | NaN | Don Mario | Takeaway/sandwich shop |
| **327** | NaN | 33-35 Moor Lane | Preston | NaN | Far East Oriental Foodstore | Retailers - other |

| | | | | | | |
|---|---|---|---|---|---|---|
| | | 393 Blackpool | | | Fredy's at | |

In [5]:

```
# We can also create a DataFrame which has only some columns
three_columns = first_1000[['RatingValue', 'FHRSID', 'PostCode']]
#열이 있는 데이터 프레임 생성

print('DataFrame with COLUMNS only:\n\n', three_columns.head())
                                                #헤드는 변수 t_c의 함수가 되며 해당 변수에 적용됨
#--> 헤드의 의미 : 처음 나오는 5개의 레코드를 보여라
    # Or some rows
# Print the Dataframe with rows 950 to 960, could be any number of rows
print("\nDataFrame with a selection of ROWS:")
first_1000[950:960]
#행 950~960 출력
```

```
DataFrame with COLUMNS only:

    RatingValue  FHRSID PostCode
0          3.0   90105  PR1 7BE
1          5.0  479725  PR1 7BY
2          5.0  479676  PR1 3BT
3          2.0  135423  PR1 6SR
4          5.0  759083      NaN

DataFrame with a selection of ROWS:
```

Out[5]:

| | AddressLine1 | AddressLine2 | AddressLine3 | AddressLine4 | BusinessName | Business |
|---|---|---|---|---|---|---|
| 950 | NaN | Sherwood Way | Preston | NaN | Sherwood Court | Hospitals Premises |
| 951 | NaN | Sherwood Way | Preston | NaN | Sherwood Lodge Res. Care Home | Hospitals Premises |
| 952 | NaN | 14 Avenham Road | Preston | NaN | Shining Stars Nursery | Hospitals Premises |
| 953 | NaN | 3 Fylde Road | Preston | NaN | Ships & Giggles | Pub/bar/ |
| 954 | NaN | 171 New Hall Lane | Preston | NaN | Shop Rite | Retailers |

| | AddressLine1 | AddressLine2 | AddressLine3 | AddressLine4 | BusinessName | Busines |
|---|---|---|---|---|---|---|
| **955** | NaN | 219 Ribbleton Lane | Preston | NaN | Sicilia Uno | Takeawa |
| **956** | NaN | NaN | NaN | NaN | Simply Country & Co | Other ca |
| **957** | NaN | NaN | NaN | NaN | Simply Delish | Other ca |
| **958** | NaN | NaN | NaN | NaN | Sion House | Hospitals Premises |
| **959** | NaN | Blackpool Road | Preston | NaN | Sir Tom Finney Community High School | School/c |

10 rows × 29 columns

Using the existing `first_1000` DataFrame, try and create a dataset which outputs the columns FHRSID, PostCode, LocalAuthorityName, with any establishment where `RatingValue < 3`

In [6]:

```
# YOUR CODE HERE
three_columns = first_1000[first_1000['RatingValue']<3][['FHRSID', 'PostCode', 'Loca
                            #뒤에 빨간색 글씨 기준을 적용해야 되기 때문에 first_1000을 다시 입력함
                            #이제 여기 뒤에 나올 FHRSID, PostCode, LocalAuthorityName 적어
print('DataFrame with COLUMNS only:\n\n', three_columns)
```

DataFrame with COLUMNS only:

```
      FHRSID PostCode LocalAuthorityName
3     135423  PR1 6SR            Preston
19    479641  PR1 2XQ            Preston
64    335295  PR1 3BT            Preston
85    448032  PR1 4DX            Preston
108   454884  PR1 1PX            Preston
120    51647  PR1 6XH            Preston
127   448082  PR1 2XQ            Preston
180   369063  PR2 6NH            Preston
189    69893  PR1 8JD            Preston
206    86973  PR1 1PX            Preston
230   637915  PR1 5LD            Preston
265   335352  PR1 8RQ            Preston
272   121600  PR1 1TS            Preston
273   133196  PR2 2DX            Preston
276   335380  PR1 7BE            Preston
292   335301  PR2 6BU            Preston
327   479591  PR1 7AT            Preston
350   201373  PR2 2DU            Preston
398   629593  PR1 5NU            Preston
426   479519  PR1 4SS            Preston
432   115804  PR1 5HH            Preston
459   479764  PR1 2AR            Preston
469   479921  PR1 4ST            Preston
474   768229  PR2 1AU            Preston
518   137209  PR1 4SU            Preston
561   850075  PR1 7JN            Preston
575   479767  PR2 1XN            Preston
632   201370  PR1 5EA            Preston
648   448112  PR1 7JS            Preston
654    70877  PR1 5QS            Preston
663    97856  PR1 8HJ            Preston
665   708327  PR1 7RA            Preston
687   335282  PR1 3BQ            Preston
700   335364  PR1 3YH            Preston
704   726619  PR1 5XA            Preston
719   479766  PR1 2EE            Preston
752   121893  PR2 6YS            Preston
787    92213  PR4 0BJ            Preston
821   136094  PR1 6EY            Preston
827   670391  PR1 6QY            Preston
828    78740  PR2 3XA            Preston
829   510354  PR1 2UP            Preston
856   479530  PR1 4TA            Preston
863   479853  PR1 5RY            Preston
868   479778  PR1 5TR            Preston
871   108298  PR1 3YH            Preston
896   136644  PR1 6QA            Preston
912    83211  PR1 5AS            Preston
```

```
925   335327   PR2 9UP           Preston
936    94388   PR1 2AR           Preston
954   479819   PR1 5XA           Preston
```

A `DataFrame` is an object in the `Pandas` library, but in addition we have the `Series` object, a collection of which makes up the `DataFrame`.

Many of the operations we can perform on a `Series` can also be performed on a `DataFrame`. It is a `Series` object which we will be using this week.

It is possible to perform an operation on each element in the `Series`, as well as call functions which require all of these such as `mean()`.

In [7]:

```python
print(first_1000['RatingValue'].head(), '\n')
print(first_1000['RatingValue'].head() * 100, '\n')
print(first_1000['RatingValue'].head() * 23 > 100)
```

```
0     3.0
1     5.0
2     5.0
3     2.0
4     5.0
Name: RatingValue, dtype: float64

0     300.0
1     500.0
2     500.0
3     200.0
4     500.0
Name: RatingValue, dtype: float64

0     False
1      True
2      True
3     False
4      True
Name: RatingValue, dtype: bool
```

As well as being a part of a `DataFrame`, it is possible to create a `Series` from a list type object, for example see the code below:

In [8]:

```
s = pd.Series([8,6,2,7,9,6]) #pd가 임포트 pandas를 의미하는거임
        #시리즈 함수를 이용하여 다음과 같은 값(리스트)을 시리즈에 할당(넣어줌)
print(type(s))
print(s)

#출력값 : 시리즈 첫번째, 멤버 0는 8이다.
```

```
<class 'pandas.core.series.Series'>
0    8
1    6
2    2
3    7
4    9
5    6
dtype: int64
```

Create a `Series` object `rating_series` which contains the `RatingValue` column from the `first_1000` `DataFrame` object.

Then display descriptive statistics from that object (mean, median, mode etc). You can see the full list of available functions in the documentation (http://pandas.pydata.org/pandas-docs/stable/generated/pandas.Series.html)

In [27]:

```
# YOUR CODE HERE
rating_series = pd.Series(first_1000['RatingValue'])
mean = rating_series.mean()
median = rating_series.median()
mode = rating_series.mode() #첫번째 최빈값은 0의 순서인 5.0이다.
                            #만약 같은 최빈값이 있다면 1의 순서인 ?가 나오겠지.

print(mean)
print(median)
print(mode)
```

```
4.29613733906
5.0
0    5.0
dtype: float64
```

# Bokeh Charts

Bokeh uses Pandas data structures as the basis for its charts. We will now place the data structure we just generated onto various types of charts.

In Bokeh, there is a class **Plot** which is the basis for the visualisations we will consider in more detail next week. A `Chart` is a subclass of `Plot`, which is designed to allow the generation of common charts with the minimum amount of code.

At its simplest, these charts simply accept an object such as a `DataFrame` or `Series`, and will output a chart. For example, a simple bar chart looks as follows:

In [2]:

```
# To display the chart in the notebook, we need to run this function, otherwise cal
output_notebook() #보케.io에서 만들어준겅
k = first_1000['RatingValue'].head()
# Create a simple bar chart with made up data
bar_chart = Bar(pd.Series(k)) #세트 혹은 리스트 멤버를 시리즈로 사용함!

# Display the bar chart
show(bar_chart) #보케.io에서 만들어준겅
```

(http://bokeh.pydata.org)

There are other charts which you can use such as a Histogram, Line graph, and scatter plot. View the Bokeh **user guide for charts (http://bokeh.pydata.org/en/latest/docs/reference/charts.html)** to see the options available to customise the chart created above. In addition it is possible to customise the display.

# Randomness

A feature which NumPy supports is that of generating random numbers. This is important, for example, in generating a random sample from a population. The randomness, however, is not *truly* random, but rather pseudo-random, i.e., it will generate predictable values based on the initial *seed* that it accepts.

This feature means that if we know the seed, we can reproduce the results we wish to share provided that we have the same data, which is a desirable property. Though this may sound counter-intuitive it allows others to run our code using the same seed and they will get the same output, the code can then be run using a different seed value.

Consider the scenario where you want to populate an array with random data, you can use the `numpy.random.randint` function as below:

In [7]:

```
# import numpy.random
import numpy as np ##넘파이에서 랜덤 만들기
# The numbers generated will include the low value
low = 0
# The numbers generated will not include the high value, but will go up to high - 1
high = 10
np.random.randint(low, high, size=10)
            #랜드int 펑션 이용하기 random의 자식임, random은 np의 자식이고 !

print(np.random.randint(low, high, size=10))
```

```
[5 6 2 3 1 7 9 9 3 2]
```

Create a loop with 10 iterations, where each iteration prints a randomly generated array of size 10. Notice how each array has set of different values.

In [8]:

```
# YOUR CODE HERE
for ri in range(0,10):
    print(np.random.randint(low, high, size=10))
```

```
[1 1 8 4 9 6 8 4 6 8]
[3 0 6 4 2 8 3 0 4 7]
[8 3 8 2 4 5 3 6 2 4]
[9 0 9 7 4 9 4 9 0 7]
[4 5 2 9 2 2 0 8 8 1]
[2 6 8 2 3 4 2 0 0 7]
[9 7 4 5 9 6 5 2 6 9]
[3 7 1 7 4 4 3 9 2 7]
[8 8 9 2 2 0 3 9 6 1]
[4 5 6 4 8 9 8 7 5 7]
```

For many situations, this is desirable. However, where we want to be able to reproduce, e.g., sample sizes, we want our samples to be reproducible. To do this, we use the `RandomState` class in NumPy, where we specify our seed, as follows:

In [9]:

```
# Run this cell several times - observe the outcome
rs = np.random.RandomState(543210)
j = rs.randint(low, high, size=10) #씨드벨류 543210을 넣으면 사용자가 전달한 시드값을 기반으로 힘
print(j)
```

```
[8 4 1 0 7 6 4 6 8 3]
```

In the cell below, generate the same loop as before, except this time instantiate the `RandomState` object to a value of 123456:

In [40]:

```
# YOUR CODE HERE
rs = np.random.RandomState(123456)
j = rs.randint(low, high, size=10)
print(j)
```

```
[1 2 1 8 0 7 4 8 4 2]
```

# IQR and Outliers

In the videos, Sergej talked about "outliers" in a dataset. In this worksheet, we'll give a slightly more detailed definition about what exactly they are, and the effect they can have on data.

An outlier is a value which is atypical of the rest of the dataset. For example, consider this set of data from searches on the UK income tax calculator (https://www.incometaxcalculator.org.uk/average-salary-uk.php). If we draw a distribution of them, we will notice a big difference in the values:

In [11]:

```python
import pandas as pd
import numpy as np
from bokeh.charts import Histogram, output_notebook, show
from bokeh.models import Axis

salaries_list =  [30000,18000,25000,20000,40000,50000,35000,45000,22000,60000,24000,
            16000,100000,21000,26000,15000,32000,19000,17000,70000,27000,55000,1850(
            36000,65000,42000,38000,12000,2481300,75000,33000,19500,43000,48000,12(
            17500,90000,34000,29000,16500,11000,31000,150000,37000,13000,22500,520(
            44000,200000,39000,46000,110000,27500,21500,47000,23500,15500,41000,265
            20500,14500,130000,250000,24500,28500,72000,140000,32500,8000,53000,95(
#50개의 결과(데이터)가 있는 검색 1번
print(salaries_list)
salaries = pd.DataFrame({'Salaries': salaries_list})

from bokeh.plotting import figure, show, output_file
output_notebook()
hist = Histogram(salaries, bins=50)
# Show absolute number on axis rather than E notation:
xaxis = hist.select(dict(type=Axis))[0] ##축 서식 지정에 관한 정보
xaxis.formatter.use_scientific = False ##축 서식 지정에 관한 정보
show(hist)
```
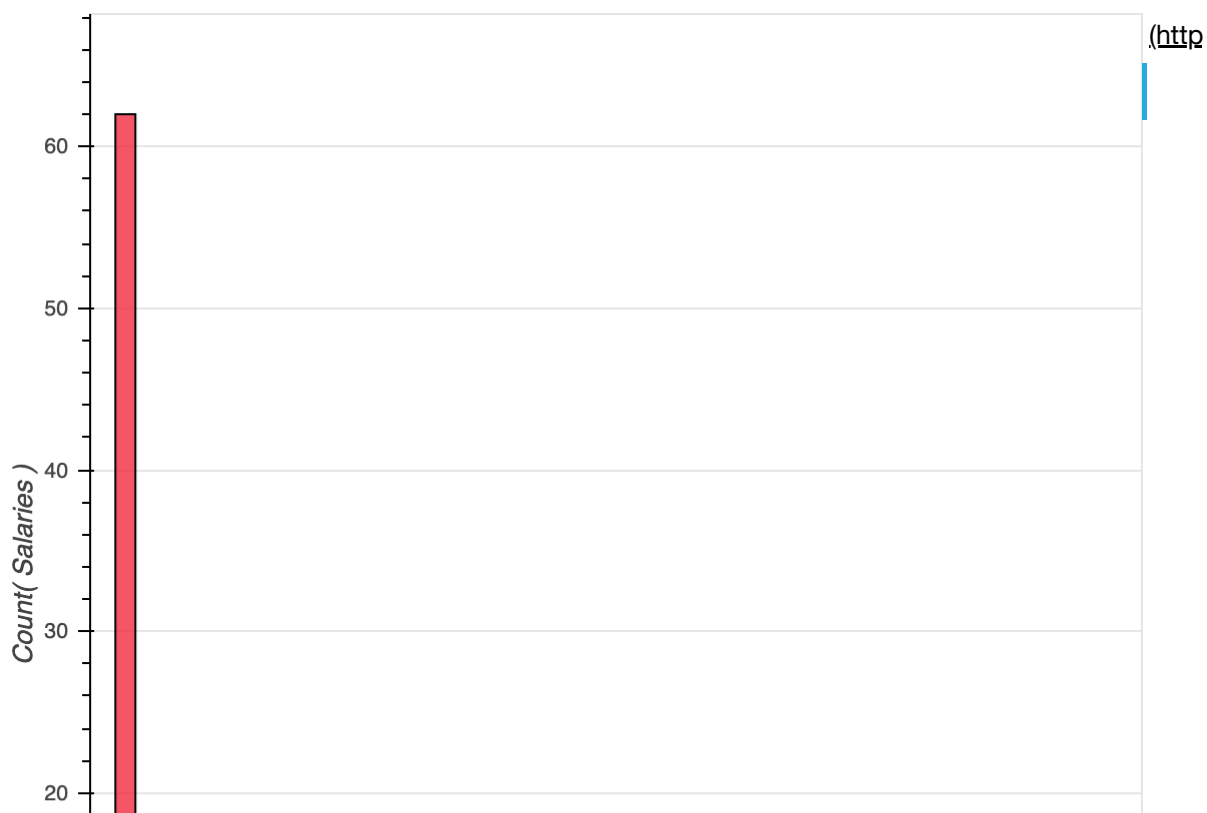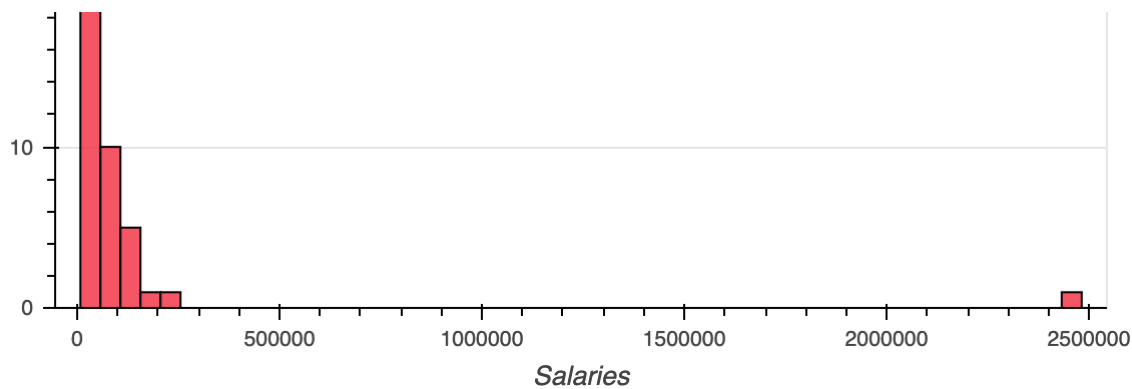
```
[30000, 18000, 25000, 20000, 40000, 50000, 35000, 45000, 22000, 60000,
24000, 28000, 23000, 16000, 100000, 21000, 26000, 15000, 32000, 19000,
17000, 70000, 27000, 55000, 18500, 80000, 36000, 65000, 42000, 38000,
12000, 2481300, 75000, 33000, 19500, 43000, 48000, 120000, 14000, 1750
0, 90000, 34000, 29000, 16500, 11000, 31000, 150000, 37000, 13000, 225
00, 52000, 10000, 85000, 44000, 200000, 39000, 46000, 110000, 27500, 2
1500, 47000, 23500, 15500, 41000, 26500, 15600, 16800, 20500, 14500, 1
30000, 250000, 24500, 28500, 72000, 140000, 32500, 8000, 53000, 95000,
25500]
```

(http:Book:beHSpyolatesafully loaded.
.(http:Book:beHSpyolatesafull

You will notice, the massive outlier on the right, where the person in question earns nearly £2.5 million. It makes it very difficult to get the chart to display anything useful, and has a significant effect on our data. For example, see the code below which shows the difference in number between the mean and the median:

In [12]:

```
print('The mean is: %f, and the median is %f' % (salaries.mean(), salaries.median())
```

The mean is: 76371.250000, and the median is 31500.000000

Task: Remove the highest value from the dataset and see how this changes the mean and the median.

Further information about explicit location based indexing .loc can be found on this Pandas Documentation (http://pandas.pydata.org/pandas-docs/version/0.15.0/indexing.html) page.

In [21]:

```
# YOUR CODE HERE 높은 값 제거
salaries = salaries.loc[salaries['Salaries'] < 2000000]
# loc[]함수는 판다스의 위치기반 인덱싱 기능
print('The mean is: %f, and the median is %f' % (salaries.mean(), salaries.median()))
```

```
-------------------------------------------------------------------
-----
TypeError                                 Traceback (most recent call
 last)
pandas/index.pyx in pandas.index.IndexEngine.get_loc (pandas/index.c:4
279)()

pandas/src/hashtable_class_helper.pxi in pandas.hashtable.Int64HashTab
le.get_item (pandas/hashtable.c:8543)()

TypeError: an integer is required

During handling of the above exception, another exception occurred:

KeyError                                  Traceback (most recent call
 last)
<ipython-input-21-e60196e84ad2> in <module>()
      1 # YOUR CODE HERE 높은 값 제거
----> 2 salaries = salaries.loc[salaries['Salaries'] < 2000000]
      3 # loc[]함수는 판다스의 위치기반 인덱싱 기능
      4 print('The mean is: %f, and the median is %f' % (salaries.mean
(), salaries.median()))
      5 print(salaries)

/opt/conda/lib/python3.5/site-packages/pandas/core/series.py in __geti
tem__(self, key)
    601           key = com._apply_if_callable(key, self)
    602           try:
--> 603               result = self.index.get_value(self, key)
    604
    605               if not is_scalar(result):

/opt/conda/lib/python3.5/site-packages/pandas/indexes/base.py in get_v
alue(self, series, key)
   2167           try:
   2168               return self._engine.get_value(s, k,
-> 2169                                    tz=getattr(series.dt
ype, 'tz', None))
   2170           except KeyError as e1:
   2171               if len(self) > 0 and self.inferred_type in ['integ
er', 'boolean']:

pandas/index.pyx in pandas.index.IndexEngine.get_value (pandas/index.
c:3557)()

pandas/index.pyx in pandas.index.IndexEngine.get_value (pandas/index.
c:3240)()

pandas/index.pyx in pandas.index.IndexEngine.get_loc (pandas/index.c:4
363)()

KeyError: 'Salaries'
```

Moving the top value had a considerable effect on the mean value of the dataset, decreasing it by over £30,000, however the result is still quite a bit higher than the median. So although the £2.5 million figure is obviously an outlier, how can we define an outlier more concretely? To start, we will consider the interquartile range (IQR):

## IQR

The IQR is calculated as follows:

1. Ordering the data by value
2. Taking the middle value from the *bottom* half of the data (lower quartile, known as Q1)
3. The median is known as Q2
4. Taking the middle value from the *top* half of the data (upper quartile, known as Q3)
5. The IQR is then calculated with Q3 - Q1

The Q1 and Q2 values are considered as the 25th and 75th percentiles, since they represent the values 25% and 75% through the ordered data. Luckily, there are functions within Pandas which allow the calculation of these percentiles, which provide us with the IQR: the `quantile` (http://pandas.pydata.org/pandas-docs/stable/generated/pandas.DataFrame.quantile.html) function on a `DataFrame` which takes a float between 0 and 1 to get the appropriate percentile. For example, the median could be calculated as follows:

In [14]:

```
salaries.quantile(0.5) #중앙값, 중앙간 사분위수 31000
```

Out[14]:

```
Salaries     31000.0
Name: 0.5, dtype: float64
```

Task: Calculate the IQR of the salaries data using Pandas. Print out the value of the upper and lower quartiles so you can check the answer is correct

In [15]:

```
# YOUR CODE HERE
upper = float(salaries.quantile(0.75)) #상위 사분위수 51000
lower = float(salaries.quantile(0.25)) #하위 사분위수 20250
iqr = upper - lower #사분범위 IQR = 30750
print('Upper quartile:\t%f\nLower quartile:\t%f\nIQR:\t\t%f' % (upper, lower, iqr))
```

```
Upper quartile: 51000.000000
Lower quartile: 20250.000000
IQR:            30750.000000
```

## Outliers

Having introduced the IQR, we can now consider what constitues an outlier. As a rule of thumb, an outlier can be defined as follows:

- `lower_quartile - (1.5 * IQR)`
- `upper_quartile + (1.5 * IQR)`

This is highly dependent on the data, and may not be appropriate for all situations, as is the decision of what to do with them. For the time being, we will simply exclude data which are outside these limits.

To do this, consider the following Pandas code, which excludes outliers from the salaries data. It uses a more complicated version of `.loc`, where it filters on two conditions

In [22]:

```python
## 2000000 제거 하는 거 보다 더 좋은 방법 !

## lower_quartile - (1.5 * IQR), upper_quartile + (1.5 * IQR) 한계치 보다
## (상한 하한 보다 ) 높거나 낮은 모든 것은 이상치로 간주함, 항상 적용되는건 아님
#You don't need to write anything here
# Create the dataset again, rather than use the one with the top value taken out
salaries = pd.DataFrame({'Salaries': salaries_list})

salaries = salaries['Salaries'][(salaries['Salaries'] > (float(lower) - (iqr * 1.5))
                    & (salaries['Salaries'] < (float(upper) + (iqr * 1.5)))]

print('The mean is: %f, and the median is %f' % (salaries.mean(), salaries.median())

print(salaries)
hist = Histogram(salaries, bins=50)
show(hist)
```

```
The mean is: 34202.816901, and the median is 28000.000000
0      30000
1      18000
2      25000
3      20000
4      40000
5      50000
6      35000
7      45000
8      22000
9      60000
10     24000
11     28000
12     23000
13     16000
15     21000
16     26000
17     15000
18     32000
19     19000
20     17000
21     70000
22     27000
23     55000
24     18500
25     80000
26     36000
27     65000
28     42000
29     38000
30     12000
       ...
44     11000
45     31000
47     37000
48     13000
49     22500
50     52000
51     10000
52     85000
53     44000
55     39000
```
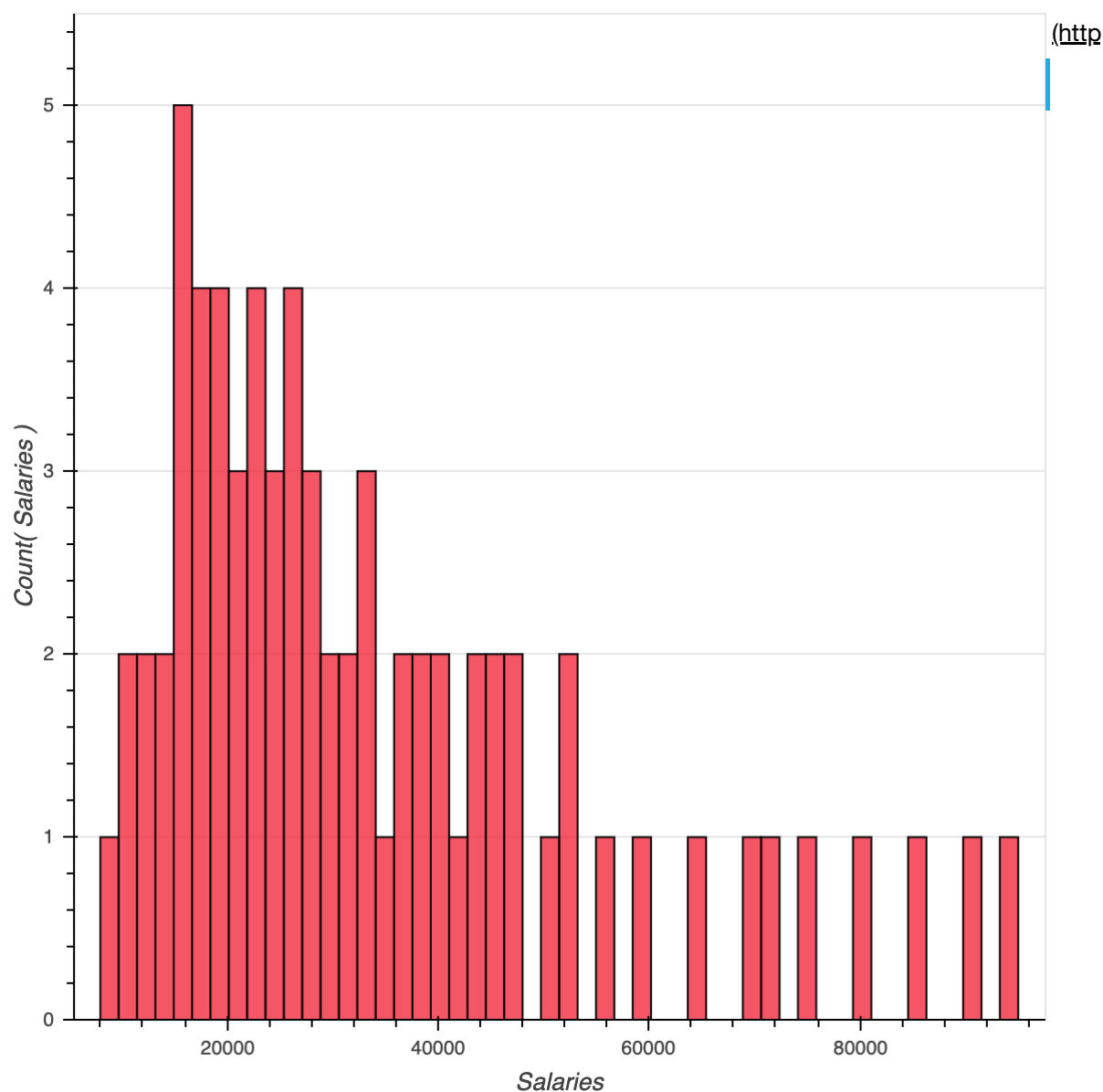
```
56    46000
58    27500
59    21500
60    47000
61    23500
62    15500
63    41000
64    26500
65    15600
66    16800
67    20500
68    14500
71    24500
72    28500
73    72000
75    32500
76     8000
77    53000
78    95000
79    25500
Name: Salaries, dtype: int64
```



The purpose of this exercise was to introduce the concept of an outlier, and how much of an effect it can have on data, and to give some practice using Pandas. There are many different ways that outliers could be defined,

and circumstances where they could or should not be excluded.