

CS156 Assignment 5: Density Estimation

Yoav Rabinovich, November 2018

For the purposes of this assignment I wrote a density estimation function that performs grid search cross validation to find the best estimator from a variety of kernels and bandwidths, and plots the data. I fed it transaction counts by month and by day of the month, as well as transaction size, to construct density estimations. I then sample a fraudulent month of transactions, by sampling a count for monthly transactions, then sampling this many day tags from our day of the month distribution, and then sampling from our size distribution for each transaction.

However, the assignment instructions left out a crucial step. As it is, we're not creating a distribution of transaction size based on the day of the month. This makes the second stage quite useless, since the information we gather from sampling days of the month isn't utilized, and we might as well have drawn transaction sizes in the amount we drew from our monthly distribution. This is a flaw that will be utilized by fraud researchers to take us down. There might be other correlations we haven't taken into consideration, for example transaction size and month of the year, since for example we buy more presents around christmas. Each such possible causal information is a tool that can be used against us, so I would refrain for using this particular model for my personal real life attempts at tax fraud.

However, our Benford analysis does show a striking proximity between our data and a true Benford distribution, which means our estimators managed to capture that feature of real life transactions, and probably others we haven't explicitly addressed.

In [431]:

```
1 # Imports
2 import numpy as np
3 import pandas as pd
4 import matplotlib.pyplot as plt
5 from scipy.stats import norm
6 from sklearn.neighbors.kde import KernelDensity
7 from sklearn.model_selection import GridSearchCV
```

In [348]:

```
1 # Read data
2 data = pd.read_csv("anonymized.csv")
3 data.Date=pd.to_datetime(data.Date,infer_datetime_format=True)
```

In [419]:

```
1
2 """
3 Plotting adapted from:
4 https://scikit-learn.org/stable/auto\_examples/neighbors/plot\_kde\_1d.html
5 """
6
7 def density(data, title=None, resolution=1000,
8             kernels=['gaussian', 'tophat', 'epanechnikov', "exponential"]):
9
10     data = data.values[:, np.newaxis]
11
12     minval = np.min(data)
13     maxval = np.max(data)
14     buffer = maxval*0.2
15
16     X_plot = np.linspace(minval-buffer, maxval+buffer, resolution)[:, np.newaxis]
17     fig, ax = plt.subplots()
18
19     maxpoint=0
20     best_est=None
21     best_score=0
22
23     for kernel in kernels:
24         grid = GridSearchCV(KernelDensity(kernel=kernel),
25                             {'bandwidth': np.logspace(-1, 5, 20)},
26                             cv=3)
27         grid.fit(data)
28         kde = grid.best_estimator_
29         log_dens = kde.score_samples(X_plot)
30         ax.plot(X_plot[:, 0], np.exp(log_dens), '-',
31                 label="kernel = '{0}'".format(kernel))
32         if max(np.exp(log_dens)) > maxpoint: maxpoint=max(np.exp(log_dens))
33         if (best_est == None) or (best_score < grid.best_score_):
34             best_est = kde
35             best_score = grid.best_score_
36
37     plt.title(title)
38     ax.legend(loc='upper right')
39     scatterpos = 0-maxpoint*0.05
40     scatterscatter= maxpoint*0.02
41     ax.plot(data[:, 0],
42             scatterpos - scatterscatter * np.random.random(data.shape[0]),
43             '.k')
44
45     ax.set_xlim(minval-buffer, maxval+buffer)
46     ax.set_ylim(scatterpos*2, maxpoint*1.1)
47     plt.show()
48     return best_est
```

In [421]:

```
1 # Prepare our datasets
2 monthly = data['Date'].groupby([data.Date.dt.year, data.Date.dt.month]).agg('count')
3 daily = data['Date'].dt.day
4 amountwise = data['Amount']
5
6 # Perform density estimation
7 monthly_estimator=density(monthly,"Q1: Density of Transaction Count per Month")
8 # Notice that the daily distribution is a simple discrete multinomial distribution,
9 # so no fancy interpolation is required.
10 daily_estimator=density(daily,"Q2: Density of Transaction Count per Day of the Month",
11                          kernels=["tophat"])
12 amountwise_estimator=density(amountwise,"Q3: Density of Transaction Amount")
13
14 #The deprecation warnings were unavoidable
15 # since they're incurred from within the GridSearchCV method
```

C:\Users\rabin\anaconda3\lib\site-packages\sklearn\model_selection_search.py:841: DeprecationWarning: The default of the `iid` parameter will change from True to False in version 0.22 and will be removed in 0.24. This will change numeric results when test-set sizes are unequal.

DeprecationWarning)

C:\Users\rabin\anaconda3\lib\site-packages\sklearn\model_selection_search.py:841: DeprecationWarning: The default of the `iid` parameter will change from True to False in version 0.22 and will be removed in 0.24. This will change numeric results when test-set sizes are unequal.

DeprecationWarning)

C:\Users\rabin\anaconda3\lib\site-packages\sklearn\model_selection_search.py:791: RuntimeWarning: invalid value encountered in subtract

array_means[:, np.newaxis]) ** 2,

C:\Users\rabin\anaconda3\lib\site-packages\sklearn\model_selection_search.py:841: DeprecationWarning: The default of the `iid` parameter will change from True to False in version 0.22 and will be removed in 0.24. This will change numeric results when test-set sizes are unequal.

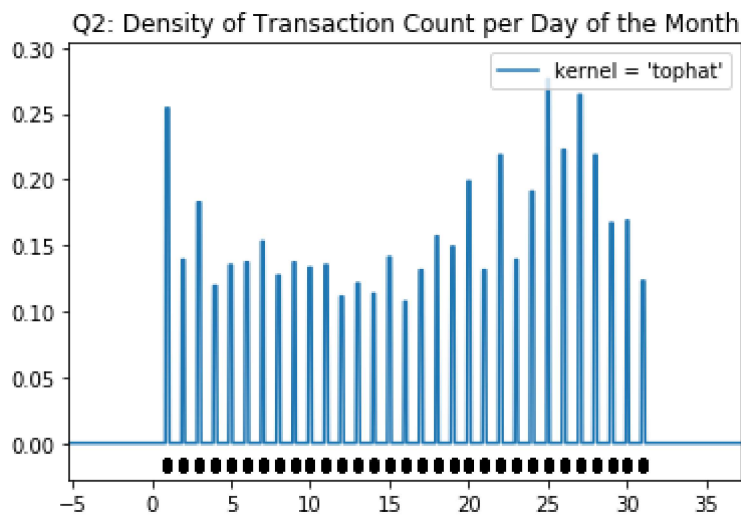
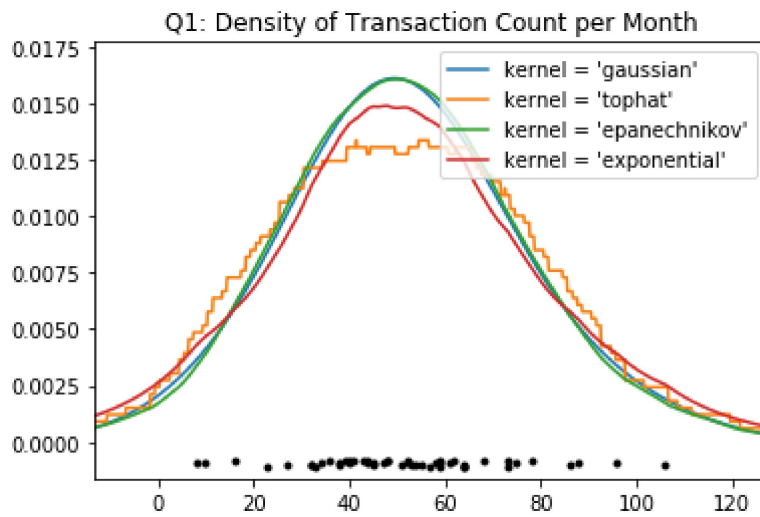
DeprecationWarning)

C:\Users\rabin\anaconda3\lib\site-packages\sklearn\model_selection_search.py:791: RuntimeWarning: invalid value encountered in subtract

array_means[:, np.newaxis]) ** 2,

C:\Users\rabin\anaconda3\lib\site-packages\sklearn\model_selection_search.py:841: DeprecationWarning: The default of the `iid` parameter will change from True to False in version 0.22 and will be removed in 0.24. This will change numeric results when test-set sizes are unequal.

DeprecationWarning)



C:\Users\rabin\anaconda3\lib\site-packages\sklearn\model_selection_search.p
y:841: DeprecationWarning: The default of the `iid` parameter will change fr
om True to False in version 0.22 and will be removed in 0.24. This will chan
ge numeric results when test-set sizes are unequal.

DeprecationWarning)

C:\Users\rabin\anaconda3\lib\site-packages\sklearn\model_selection_search.p
y:791: RuntimeWarning: invalid value encountered in subtract

array_means[:, np.newaxis]) ** 2,

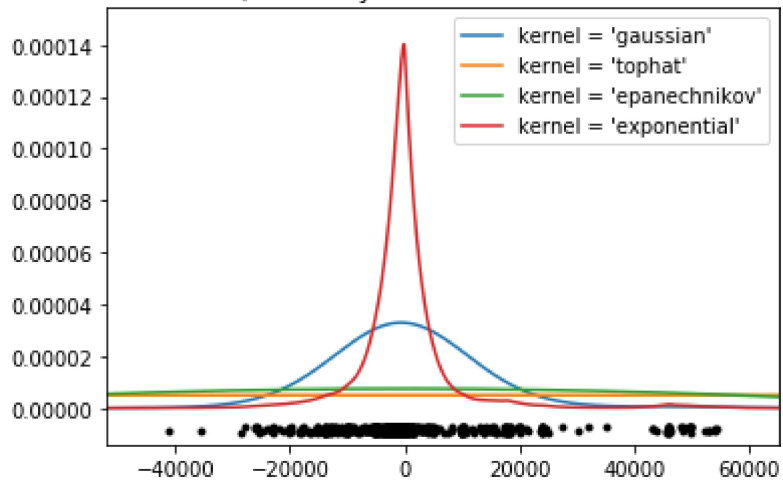
C:\Users\rabin\anaconda3\lib\site-packages\sklearn\model_selection_search.p
y:791: RuntimeWarning: invalid value encountered in subtract

array_means[:, np.newaxis]) ** 2,

C:\Users\rabin\anaconda3\lib\site-packages\sklearn\model_selection_search.p
y:841: DeprecationWarning: The default of the `iid` parameter will change fr
om True to False in version 0.22 and will be removed in 0.24. This will chan
ge numeric results when test-set sizes are unequal.

DeprecationWarning)

Q3: Density of Transaction Amount



In [422]:

```
1 # Producing a fraudulent sample month of transactions
2 # Forcing Gaussian and Tophat,
3 # since other kernels aren't supported by the sampling function
4 monthly_estimator=density(monthly,kernels=["gaussian","tophat"]);
5 daily_estimator=density(daily,kernels=["tophat"]);
6 amountwise_estimator=density(amountwise,kernels=["gaussian","tophat"]);
```

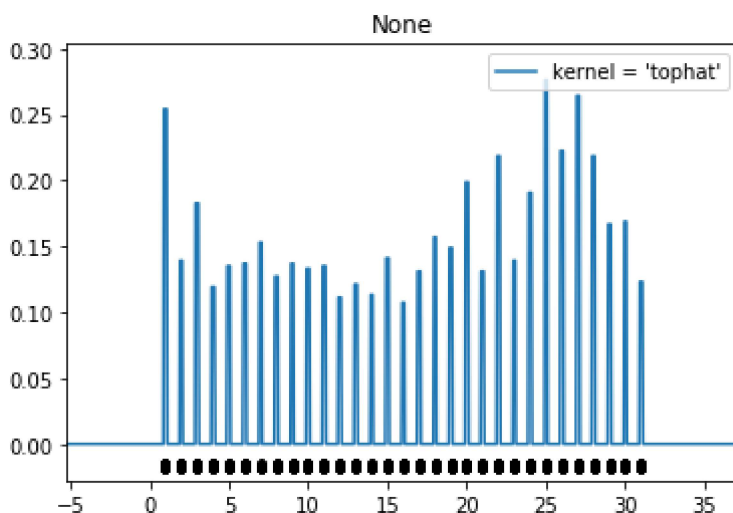
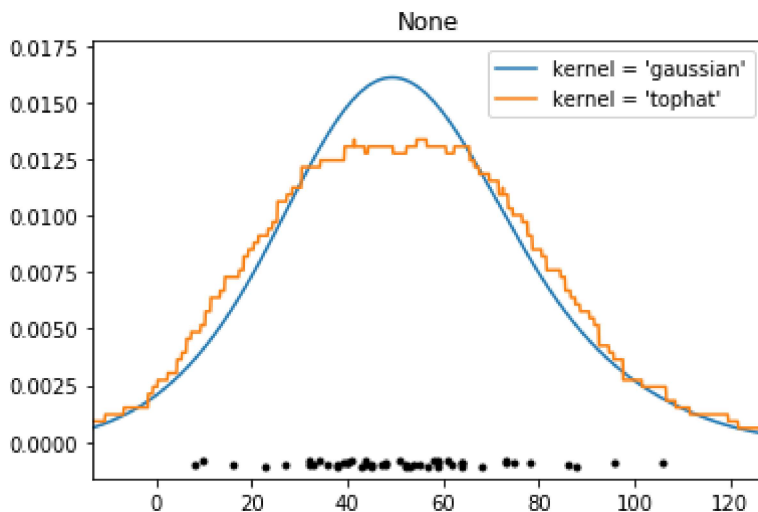
C:\Users\rabin\anaconda3\lib\site-packages\sklearn\model_selection_search.p
y:841: DeprecationWarning: The default of the `iid` parameter will change fr
om True to False in version 0.22 and will be removed in 0.24. This will chan
ge numeric results when test-set sizes are unequal.

DeprecationWarning)

C:\Users\rabin\anaconda3\lib\site-packages\sklearn\model_selection_search.p
y:841: DeprecationWarning: The default of the `iid` parameter will change fr
om True to False in version 0.22 and will be removed in 0.24. This will chan
ge numeric results when test-set sizes are unequal.

DeprecationWarning)

C:\Users\rabin\anaconda3\lib\site-packages\sklearn\model_selection_search.p
y:791: RuntimeWarning: invalid value encountered in subtract
array_means[:, np.newaxis]) ** 2,



C:\Users\rabin\anaconda3\lib\site-packages\sklearn\model_selection_search.p

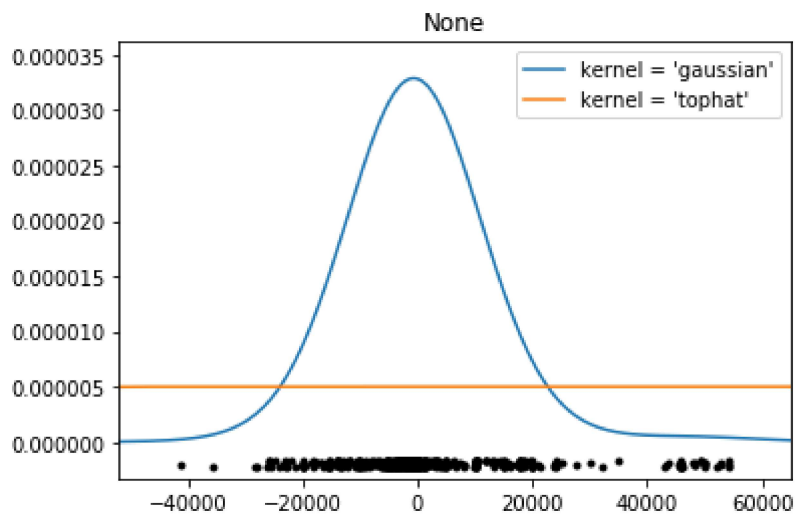
y:841: DeprecationWarning: The default of the `iid` parameter will change from True to False in version 0.22 and will be removed in 0.24. This will change numeric results when test-set sizes are unequal.

DeprecationWarning)

C:\Users\rabin\anaconda3\lib\site-packages\sklearn\model_selection_search.p

y:791: RuntimeWarning: invalid value encountered in subtract

array_means[:, np.newaxis]) ** 2,



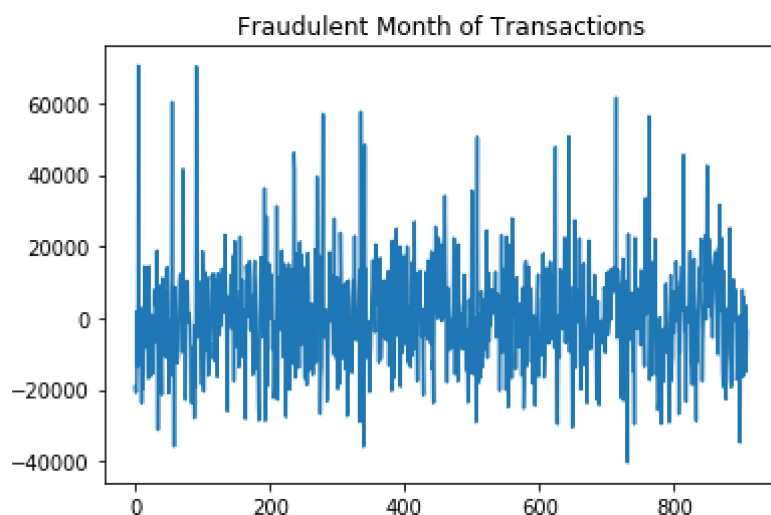
»

In [485]:

```
1 # Drawing a monthly count
2 new_monthly = np.round(monthly_estimator.sample(1)).astype(int)
3 # Drawing day tags
4 new_daily = np.round(daily_estimator.sample(new_monthly[0][0])).astype(int)
5 # Drawing transaction sizes
6 new_amounts = np.around(amountwise_estimator.sample(sum(new_daily)), decimals=2)
7
8 plt.title("Fraudulent Month of Transactions")
9 plt.plot(new_amounts)
```

Out[485]:

[<matplotlib.lines.Line2D at 0x26b06e91da0>]



In [499]:

```
1 # Checking for conformity with Benford's Law
2
3 benford = np.zeros(10)
4
5 for trans in new_amounts:
6     first_digit=int(str(abs(trans))[1])
7     benford[first_digit-1]+=1
8
9 index = range(1,11)
10 true = [np.log10(1 + (1.0 / d)) for d in index]
11 plt.subplot(1,2,1)
12 plt.bar(index,benford,color='b')
13 plt.subplot(1,2,2)
14 plt.bar(index,true,color='r')
15 plt.show()
```

