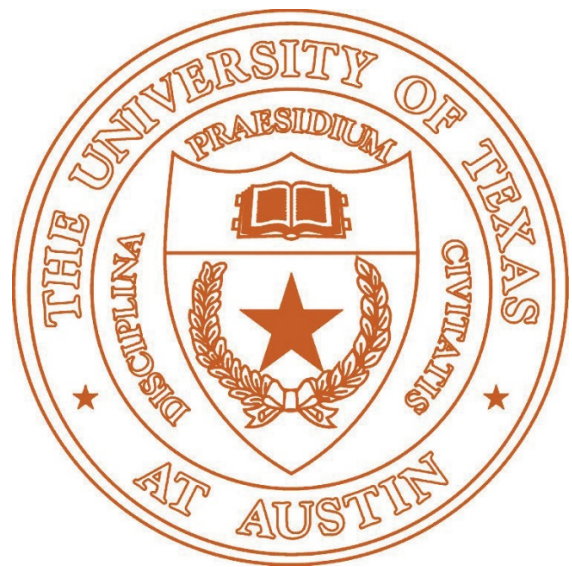


University of Texas at Austin, Cockrell School of Engineering
Data Mining – EE 380L



Problem Set # 3

April 11, 2016

Gabrielson Eapen

EID: EAPENGP

Discussed Homework with Following Students:

1. Mudra Gandhi
2. Rayo Landeros

Q1]

```
In [1]: # Name: Gabe Eapen
# UT EID: eapengp
# PS3 - Q1

In [2]: import pandas as pd
import numpy as np
import matplotlib.pyplot as plt
%matplotlib inline
from sklearn import datasets, linear_model
from pandas import DataFrame, Series
import seaborn as sns
sns.set(style='ticks', palette='Set2')

In [3]: def extract_int(some_string):
    int_as_string = (str(some_string)).split('.')[0]
    return int(int_as_string)

In [4]: df=pd.read_stata("nes5200_processed_voters_realideo.dta")
df.shape

Out[4]: (41498, 62)

In [5]: print(df.columns)

Index([u'year', u'resid', u'weight1', u'weight2', u'weight3', u'age',
      u'gender', u'race', u'educ1', u'urban', u'region', u'income', u'occup1',
      u'union', u'religion', u'educ2', u'educ3', u'martial_status', u'occup2',
      u'icper_cty', u'fips_cty', u'partyid7', u'partyid3', u'partyid3_b',
      u'str_partyid', u'father_party', u'mother_party', u'dlikes', u'rlikes',
      u'dem_therm', u'rep_therm', u'regis', u'vote', u'regisvote',
      u'presvote', u'presvote_2party', u'presvote_intent', u'ideo_feel',
      u'ideo7', u'ideo', u'cd', u'state', u'inter_pre', u'inter_post',
      u'black', u'female', u'age_sq', u'rep_presvote', u'rep_pres_intent',
      u'south', u'real_ideo', u'presapprov', u'perfin1', u'perfin2',
      u'perfin', u'presadm', u'age_10', u'age_sq_10', u'newfathe', u'newmoth',
      u'parent_party', u'white'],
      dtype='object')

In [6]: df_1992_raw = df[(df['year'] == 1992.0) & ((df['presvote'] == "1. democrat") | (df['presvote'] =
print df_1992_raw.shape
df_ed1992 = df.loc[(df['year'] == 1992.0) & ((df['presvote'] == "1. democrat") | (df['presvote']
print df_ed1992.shape

(1304, 62)
(1304, 62)
```

Part a)

```
In [7]: df_vote_inc = pd.DataFrame(df_ed1992, columns=['presvote', 'income'])
print df_vote_inc.shape
print df_vote_inc.head()
df_clean = df_vote_inc.dropna(how='any')
print df_clean.shape
#print df_clean.dtypes
print df_clean.head()
```

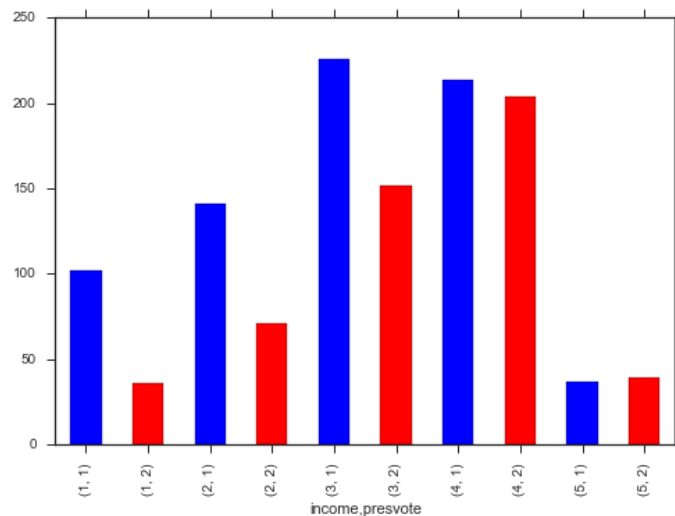
```
(1304, 2)
      presvote      income
32092  2. republican  4. 68 to 95 percentile
32093  2. republican  2. 17 to 33 percentile
32095  1. democrat   1. 0 to 16 percentile
32096  2. republican  2. 17 to 33 percentile
32097  1. democrat   3. 34 to 67 percentile
(1222, 2)
presvote    category
income      category
dtype: object
      presvote      income
32092  2. republican  4. 68 to 95 percentile
32093  2. republican  2. 17 to 33 percentile
32095  1. democrat   1. 0 to 16 percentile
32096  2. republican  2. 17 to 33 percentile
32097  1. democrat   3. 34 to 67 percentile
```

```
In [8]: cat_columns = df_clean.select_dtypes(['category']).columns
#cat_columns
```

```
In [9]: df_clean[cat_columns] = df_clean[cat_columns].apply(lambda x: x.cat.codes + 1)
#print vote_D.head()
#df_clean.head()
```

```
In [10]: df_clean.groupby(['income', 'presvote']).size().plot(kind='bar', color=['blue', 'red'])
```

```
Out[10]: <matplotlib.axes._subplots.AxesSubplot at 0xa1db358>
```



```
Out[13]: array([2, 2, 1, ..., 2, 1, 2], dtype=int8)
```

```
Coeff: [[ 0.29072854]]
Intercept (B0) [-1.25875036]
```

[illegible]

```
In [64]: print logreg.predict_proba(XRepInc_3.reshape(-1,1))
```

[illegible]