

▼ HOMEWORK 2

```
1 import csv
2 import math as m
3 import operator
4 from scipy.stats import pearsonr
5 import matplotlib.pyplot as plt
6 import numpy as np
7 import pandas as pd
8 import matplotlib
9 from sklearn.metrics.pairwise import cosine_similarity
10 from scipy.spatial import distance
11 from scipy import spatial
12 matplotlib.rcParams.update({'font.size': 12})
13 from sklearn.datasets import load_boston
14 from sklearn.model_selection import train_test_split
15 from sklearn.linear_model import LinearRegression
16 from sklearn.linear_model import Ridge
17 from scipy.stats import pearsonr
18 from numpy import cov
19 #Creating a class for representing the Datapoints.
20
21 class Datapoint(object):
22
23     def __init__(self, feats):
24         self.feature_1 = feats['Hour']
25         self.feature_2 = feats['ImmobilizedBus']
26         self.feature_3 = feats['BrokenTruck']
27         self.feature_4 = feats['VehicleExcess']
28         self.feature_5 = feats['AccidentVictim']
29         self.feature_6 = feats['RunningOver']
30         self.feature_7 = feats['FireVehicles']
31         self.feature_8 = feats['OccurenceFreight']
32         self.feature_9 = feats['DangerousFreight']
33         self.feature_10 = feats['Electricity']
34         self.feature_11 = feats['Fire']
35         self.feature_12 = feats['PointFlooding']
36         self.feature_13 = feats['Manifestations']
37         self.feature_14 = feats['TrolleyBusNetworkERR']
38         self.feature_15 = feats['Tree']
39         self.feature_16 = feats['SemaphoreOff']
40         self.feature_17 = feats['IntermittentSemaphore']
41         self.outcome = feats['Slowness']
42     #Returns the features as a Numpy Array.
43     def feature_vector(self):
44         return np.array([self.feature_1, self.feature_2, self.feature_3,\
45                         self.feature_4, self.feature_5, self.feature_6,\
46                         self.feature_7, self.feature_8, self.feature_9, \
47                         self.feature_10,self.feature_11,self.feature_12,\
48                         self.feature_13,self.feature_14,self.feature_15,\
49                         self.feature_16,self.feature_17,self.outcome])
50
51     def __str__(self):
52         return "\Hour:{}, \nImmobilized bus:{}, \nBroken Truck:{}, \nVehicle Excess:{},\
53         \nAccident Victim:{}, \nRunning Over:{}, \nFire Vehicles:{}, \
54         \nOccurence Involving Freight:{}, \nIncident Involving Dangerous Freight:{},\
55         \nLack of Electricity:{}, \nFire:{}, \nPoint of flooding:{}, \nManifestations:{},\
56         \nDefect in the network of trolleybuses:{}, \nTree on the road:{}, \nSemaphore off:{}\
57         , \nIntermittent Semaphore:{}, \nSlowness in traffic:{} ".format(self.feature_1, self.feature_2, self.feature_3,\
58                                 self.feature_4, self.feature_5, self.feature_6,\
59                                 self.feature_7, self.feature_8, self.feature_9, \
60                                 self.feature_10,self.feature_11,self.feature_12,\
61                                 self.feature_13,self.feature_14,self.feature_15,\
62                                 self.feature_16,self.feature_17,self.outcome)
63 #Function that creates the datapoints and writes the corresponding feature value to the datapoints.
64 def parse_dataset(filename):
65     with open(filename) as csvfile:
66         dataset = []
67         lineCount = 0
68         readCSV = csv.reader(csvfile)
69         for row in readCSV:
70             if(lineCount==0):
71                 a=Datapoint({'Hour':(row[0]), 'ImmobilizedBus':(row[1]), 'BrokenTruck':(row[2]), 'VehicleExcess':(row[3]), \
72                             'AccidentVictim':(row[4]), 'RunningOver':(row[5]), 'FireVehicles':(row[6]), 'OccurenceFreight':(row[7]),\
73                             'DangerousFreight':(row[8]), 'Electricity':(row[9]), 'Fire':(row[10]), 'PointFlooding':(row[11]),\
74                             'Manifestations':(row[12]), 'TrolleyBusNetworkERR':(row[13]), 'Tree':(row[14]), 'SemaphoreOff':(row[15]), \
75                             'IntermittentSemaphore':(row[16]), 'Slowness':(row[17])})
76                 feature_list.append(a.feature_vector())
77                 lineCount += 1
78             else:
79                 a=Datapoint({'Hour':(row[0]), 'ImmobilizedBus':(row[1]), 'BrokenTruck':(row[2]), 'VehicleExcess':(row[3]), \
80                             'AccidentVictim':(row[4]), 'RunningOver':(row[5]), 'FireVehicles':(row[6]), 'OccurenceFreight':(row[7]),\
81                             'DangerousFreight':(row[8]), 'Electricity':(row[9]), 'Fire':(row[10]), 'PointFlooding':(row[11]),\
82                             'Manifestations':(row[12]), 'TrolleyBusNetworkERR':(row[13]), 'Tree':(row[14]), 'SemaphoreOff':(row[15]), \
```

```

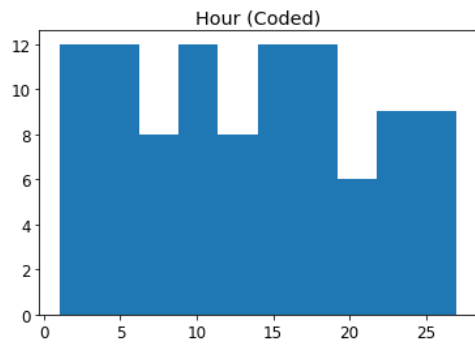
83         'IntermittentSemaphore':(row[16]),'Slowness':(row[17]))
84
85         dataset.append(a.feature_vector())
86
87
88     return dataset
89
90 feature_list = []
91 #PARSE the training , development and testing dataset.
92 dataset_train = parse_dataset('hw2__question1_train.csv')
93 dataset_test = parse_dataset('hw2__question1_test.csv')
94
95 def convertInt (dataset):
96     result = []
97     for data in dataset:
98         result.append(data.astype(int))
99     return result
100
101 dataset_test = convertInt(dataset_test)
102 dataset_train = convertInt(dataset_train)
103
104 #Printing some statistics about the data.
105 print("Total Number of Data Points in Training set: {}".format(len(dataset_train)))
106 print("Total Number of Data Points in Testing set: {}".format(len(dataset_test)))
107
108 #print(dataset_train)
109 #Function to plot.
110 def plot_histogram(dataset):
111     i = 0
112     for i in range (18):
113         print ("-----\n\
114 ",i+1,')',feature_list[0][i],"\n-----")
115         plot_list = []
116         for data in dataset:
117             plot_list.append(data[i])
118         plt.title(feature_list[0][i])
119         plt.hist(plot_list)
120         plt.show()
121
122 plot_histogram(dataset_train)
123

```

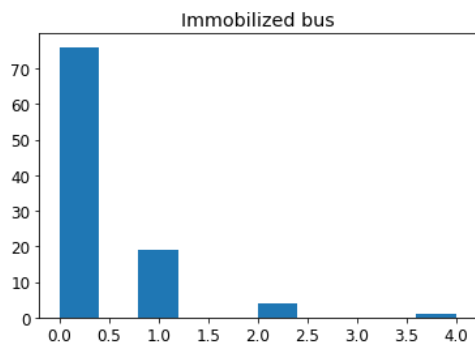


Total Number of Data Points in Training set: 100
Total Number of Data Points in Testing set: 35

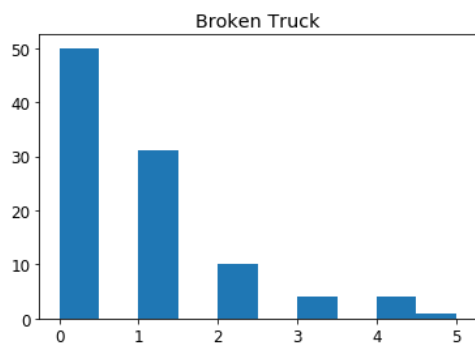
1) Hour (Coded)



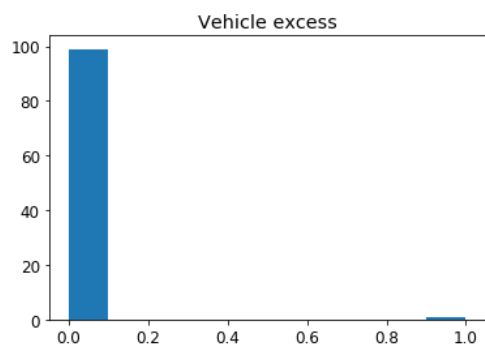
2) Immobilized bus



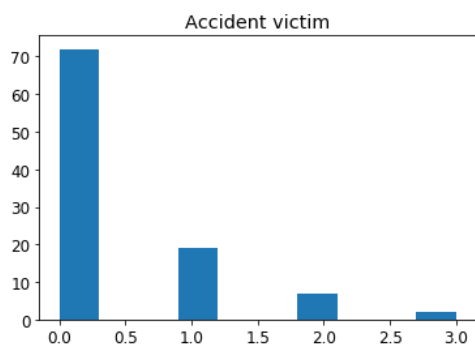
3) Broken Truck



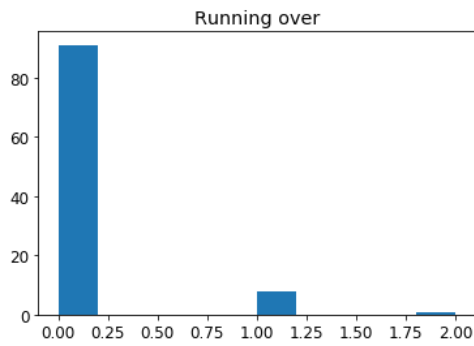
4) Vehicle excess



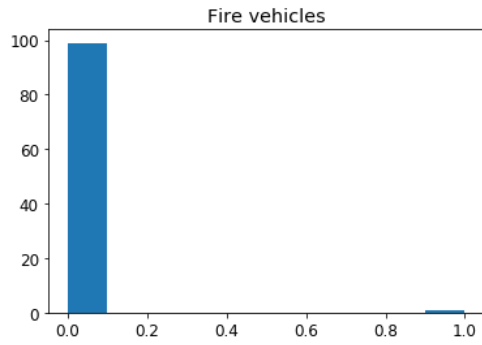
5) Accident victim



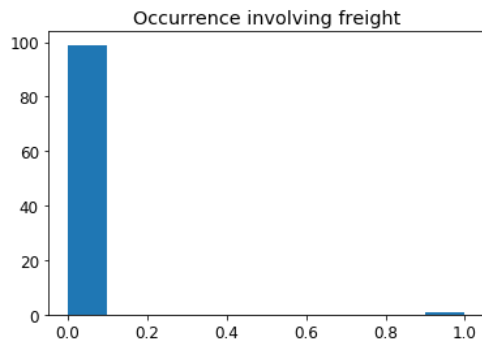
6) Running over



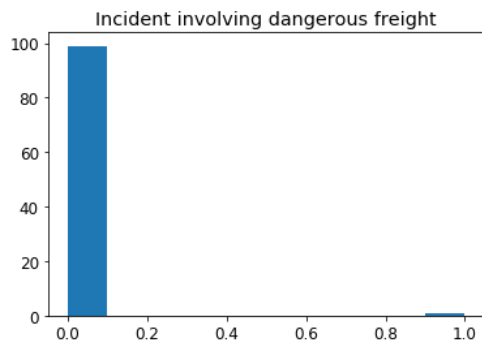
7) Fire vehicles



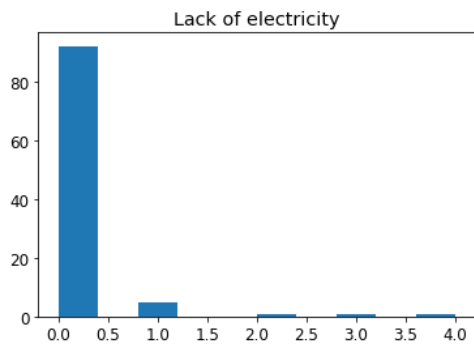
8) Occurrence involving freight



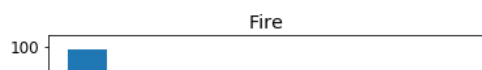
9) Incident involving dangerous freight

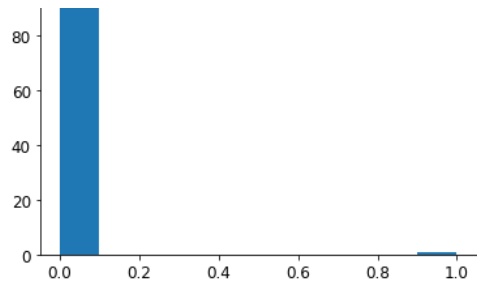


10) Lack of electricity

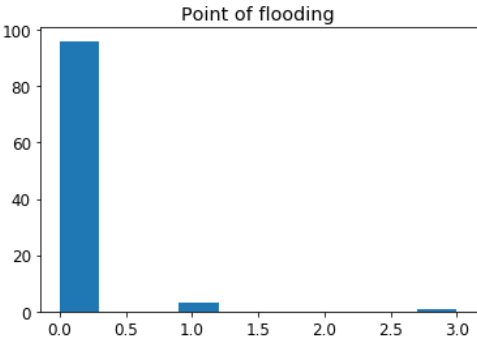


11) Fire

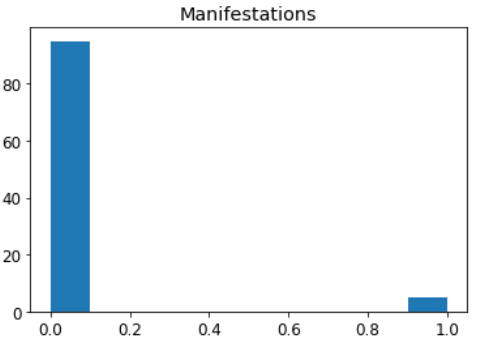




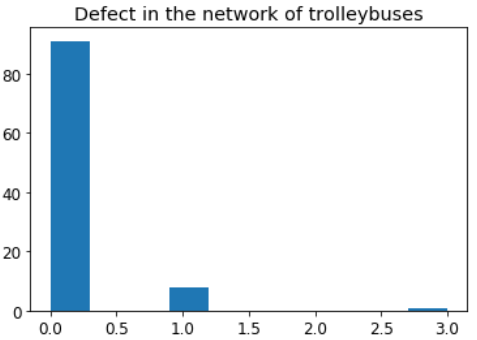
12) Point of flooding



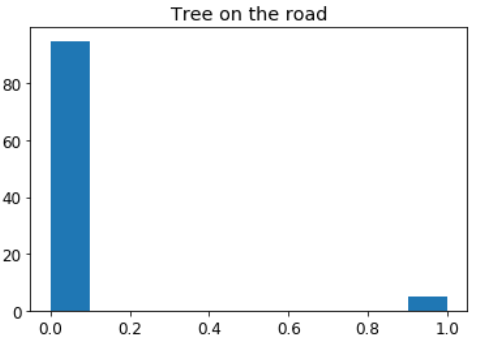
13) Manifestations



14) Defect in the network of trolleybuses

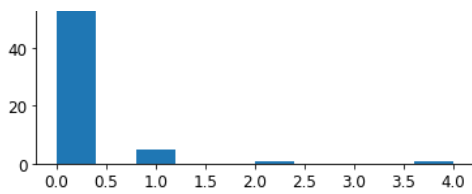


15) Tree on the road

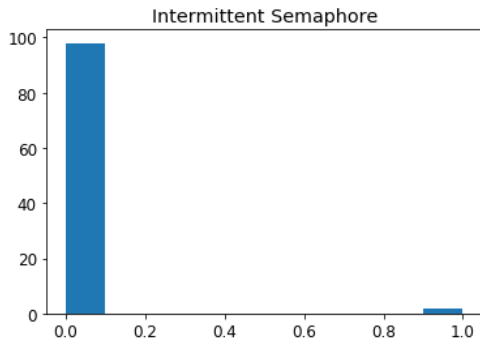


16) Semaphore off

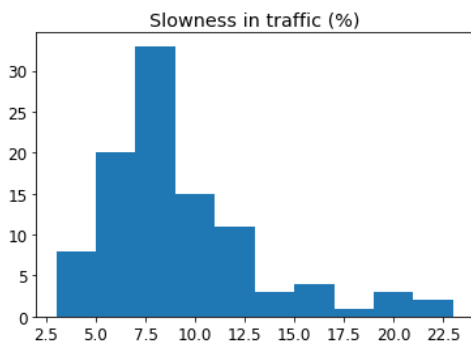




17) Intermittent Semaphore



18) Slowness in traffic (%)



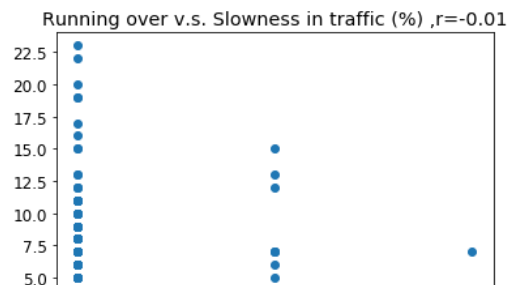
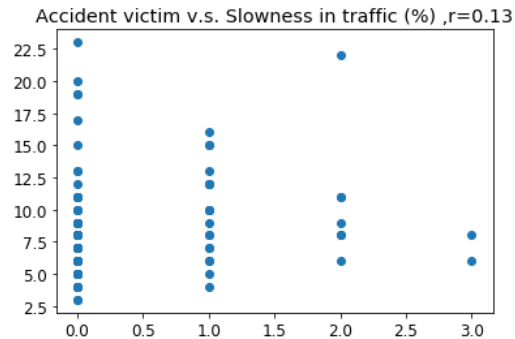
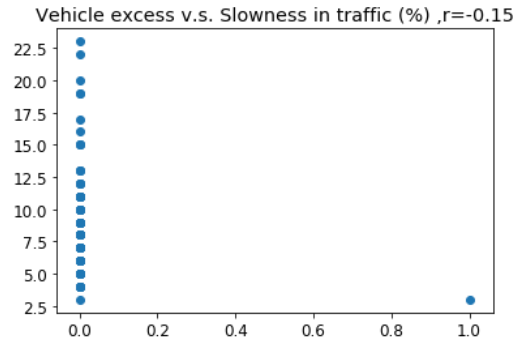
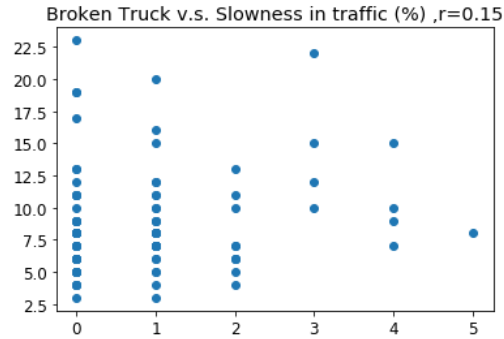
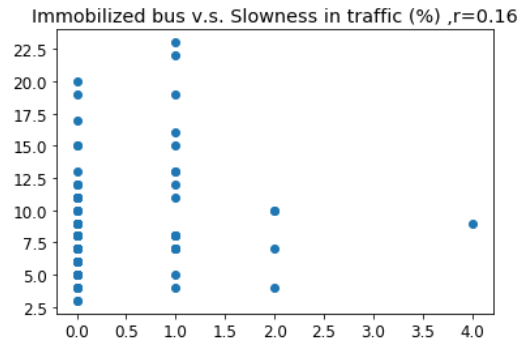
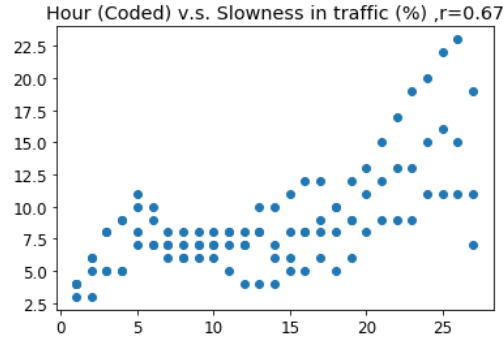
```

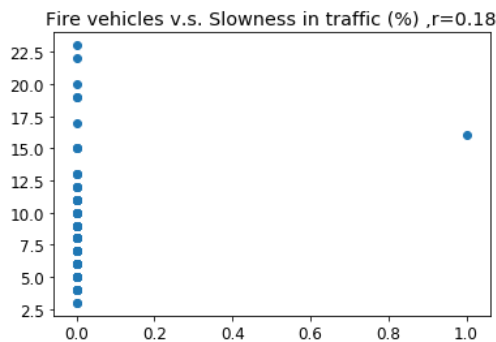
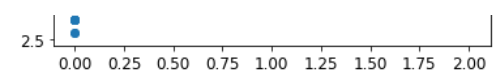
1 import pandas as pd
2
3 train_data = pd.read_csv("hw2__question1_train.csv")
4 test_data = pd.read_csv("hw2__question1_test.csv")
5
6 newX = train_data.drop('Slowness in traffic (%)',axis= 1)
7 newY = train_data['Slowness in traffic (%)']
8
9 newX_test = test_data.drop('Slowness in traffic (%)', axis= 1)
10
11 newY_test = test_data['Slowness in traffic (%)']
12
13 features_data = np.array(feature_list[0])
14 print (features_data)
15 dataset_test = np.array(dataset_test)
16 dataset_train = np.array(dataset_train)
17 trainT = dataset_train.T
18
19 for i in features_data:
20     corr, p = pearsonr(train_data[i],newY)
21     plt.scatter(train_data[i],newY)
22     plt.title(i+' v.s.'+' Slowness in traffic (%) ,r='+str(format(corr,'.2f')))
23     plt.show()
24
25 for i in features_data:
26     for j in features_data:
27         corr, p = pearsonr(train_data[i],train_data[j])
28         print(format(corr, '.2f'), end = ',')
29     print()

```

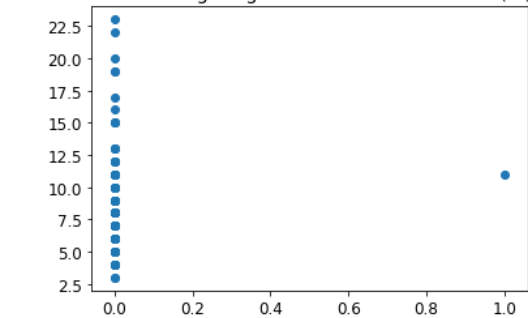


['Hour (Coded)' 'Immobilized bus' 'Broken Truck' 'Vehicle excess'
'Accident victim' 'Running over' 'Fire vehicles'
'Occurrence involving freight' 'Incident involving dangerous freight'
'Lack of electricity' 'Fire' 'Point of flooding' 'Manifestations'
'Defect in the network of trolleybuses' 'Tree on the road'
'Semaphore off' 'Intermittent Semaphore' 'Slowness in traffic (%)']

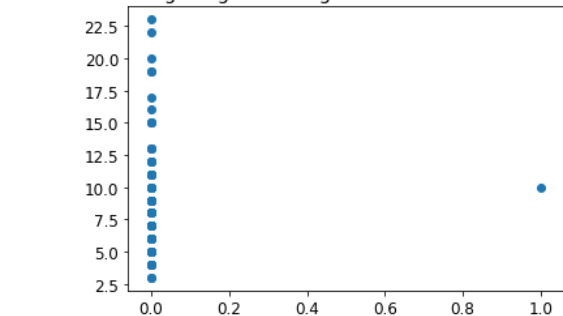




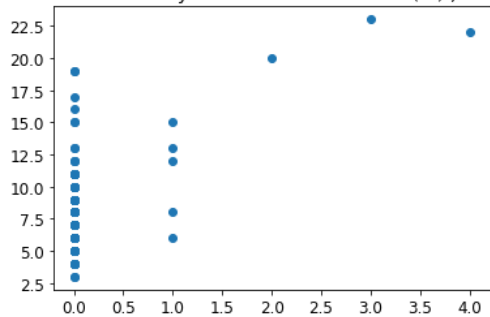
Occurrence involving freight v.s. Slowness in traffic (%) , $r=0.06$



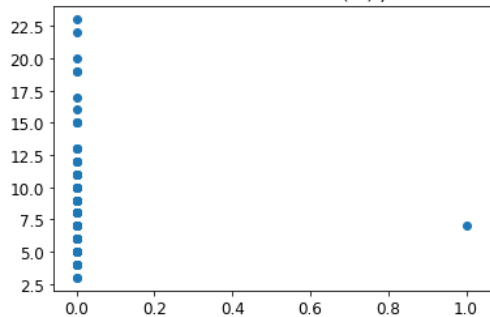
Incident involving dangerous freight v.s. Slowness in traffic (%) , $r=0.03$



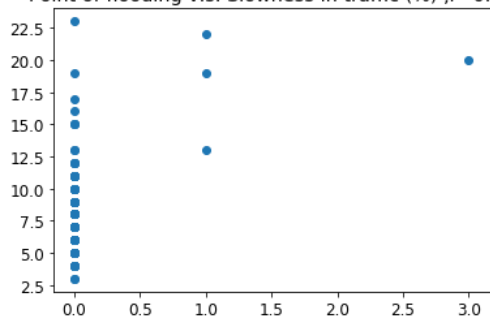
Lack of electricity v.s. Slowness in traffic (%) , $r=0.57$



Fire v.s. Slowness in traffic (%) , $r=-0.04$

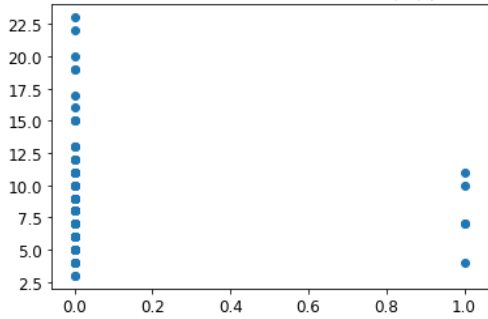


Point of flooding v.s. Slowness in traffic (%) , $r=0.46$

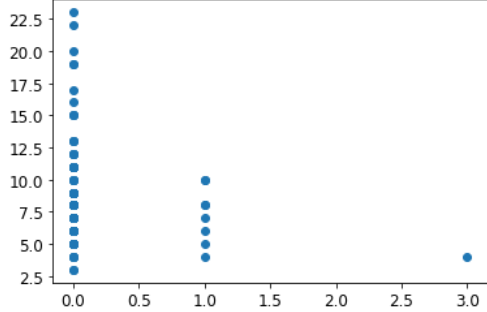


Point of flooding v.s. Slowness in traffic (%) , $r=0.46$

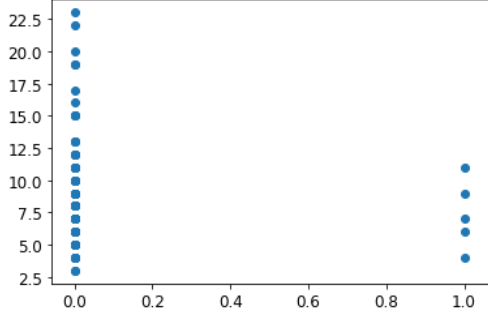
Manifestations v.s. Slowness in traffic (%) , $r=-0.06$



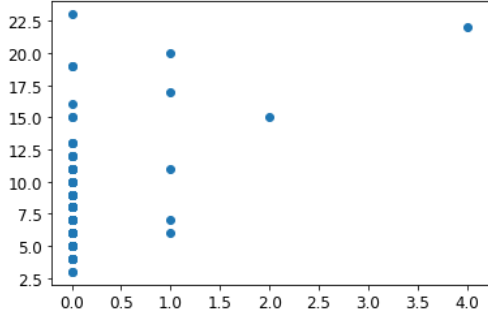
Defect in the network of trolleybuses v.s. Slowness in traffic (%) , $r=-0.17$



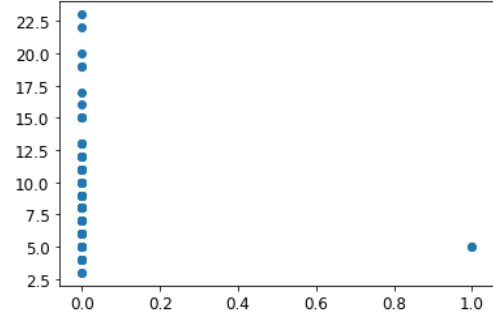
Tree on the road v.s. Slowness in traffic (%) , $r=-0.08$



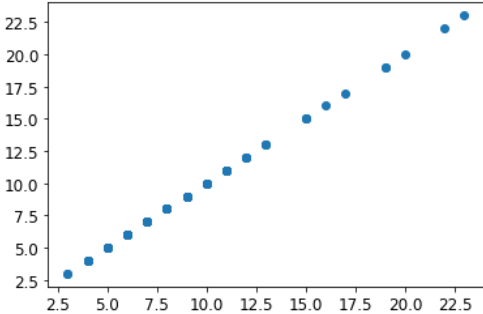
Semaphore off v.s. Slowness in traffic (%) , $r=0.43$



Intermittent Semaphore v.s. Slowness in traffic (%) , $r=-0.14$



Slowness in traffic (%) v.s. Slowness in traffic (%) , $r=1.00$



1.00,0.09,0.22,-0.16,0.19,0.04,0.16,0.02,-0.00,0.28,-0.04,0.23,-0.15,-0.18,-0.05,0.19,-0.18,0.67,
0.09,1.00,0.22,-0.05,-0.02,0.18,0.11,0.11,0.26,0.10,-0.05,0.01,0.17,0.10,0.03,0.05,0.04,0.16,
0.22,0.22,1.00,-0.08,0.40,0.15,0.01,0.10,0.10,0.11,0.01,0.08,0.03,0.11,-0.05,0.20,-0.04,0.15,
-0.16,-0.05,-0.08,1.00,-0.06,-0.03,-0.01,-0.01,-0.01,-0.02,-0.01,-0.02,-0.02,-0.03,-0.02,-0.02,-0.01,-0.15,

```

0.19,-0.02,0.40,-0.06,1.00,0.00,0.09,-0.06,0.09,0.11,-0.06,0.03,-0.06,0.06,-0.06,0.11,-0.08,0.13,
0.04,0.18,0.15,-0.03,0.00,1.00,-0.03,-0.03,-0.03,-0.02,-0.03,0.04,0.07,0.07,0.07,-0.07,0.17,-0.01,
0.16,0.11,0.01,-0.01,0.09,-0.03,1.00,-0.01,-0.01,-0.02,-0.01,-0.02,-0.02,-0.03,-0.02,-0.02,-0.01,0.18,
0.02,0.11,0.10,-0.01,-0.06,-0.03,-0.01,1.00,-0.01,-0.02,-0.01,-0.02,0.44,-0.03,-0.02,0.18,-0.01,0.06,
-0.00,0.26,0.10,-0.01,0.09,-0.03,-0.01,-0.01,1.00,-0.02,-0.01,-0.02,0.44,-0.03,-0.02,-0.02,-0.01,0.03,
0.28,0.10,0.11,-0.02,0.11,-0.02,-0.02,-0.02,1.00,-0.02,0.53,-0.06,-0.07,-0.06,0.67,-0.04,0.57,
-0.04,-0.05,0.01,-0.01,-0.06,-0.03,-0.01,-0.01,-0.02,1.00,-0.02,-0.02,-0.03,-0.02,-0.02,-0.01,-0.04,
0.23,0.01,0.08,-0.02,0.03,0.04,-0.02,-0.02,-0.02,0.53,-0.02,1.00,-0.04,-0.05,-0.04,0.38,-0.03,0.46,
-0.15,0.17,0.03,-0.02,-0.06,0.07,-0.02,0.44,0.44,-0.06,-0.02,-0.04,1.00,0.28,0.16,0.04,-0.03,-0.06,
-0.18,0.10,0.11,-0.03,0.06,0.07,-0.03,-0.03,-0.03,-0.07,-0.03,-0.05,0.28,1.00,0.05,-0.06,-0.04,-0.17,
-0.05,0.03,-0.05,-0.02,-0.06,0.07,-0.02,-0.02,-0.02,-0.06,-0.02,-0.04,0.16,0.05,1.00,-0.05,-0.03,-0.08,
0.19,0.05,0.20,-0.02,0.11,-0.07,-0.02,0.18,-0.02,0.67,-0.02,0.38,0.04,-0.06,-0.05,1.00,-0.03,0.43,
-0.18,0.04,-0.04,-0.01,-0.08,0.17,-0.01,-0.01,-0.01,-0.04,-0.01,-0.03,-0.03,-0.04,-0.03,-0.03,1.00,-0.14,
0.67,0.16,0.15,-0.15,0.13,-0.01,0.18,0.06,0.03,0.57,-0.04,0.46,-0.06,-0.17,-0.08,0.43,-0.14,1.00,

```

```

1 #Creating an LR model using OLS Solution
2 dMtx = np.array(newX)
3 dMtx=np.append(np.ones([100,1]),dMtx,axis=1)
4 dMtx = dMtx.astype('int64')
5
6 Outcome = np.array(newY)
7
8
9 first = (dMtx.T.dot(dMtx))
10 first = np.linalg.inv(first)
11 second = (dMtx.T.dot(newY))
12 w_opt = first.dot(second)
13 print(w_opt) #this is the OLS model.
14
15 #Now using this model on our TEST DATA and calculate Pearson's r and RSS.
16
17 dMtx_test = np.array(newX_test)
18 dMtx_test = np.append(np.ones([35,1]),dMtx_test,axis=1)
19 dMtx_test = dMtx_test.astype('int64')
20
21 test_outcome = np.array(newY_test)
22 RSS = 0
23
24 inner = dMtx_test.dot(w_opt)
25
26 first = (test_outcome-inner)
27 RSS = first.T.dot(first)
28
29 print('RSS ERROR=',RSS)
30
31 test_score, test_score_p = pearsonr(inner,test_outcome)
32 print("r=" + format(test_score, '.2f') + ", p=" + format(test_score_p, '.2f'))
33
34

```

```

[ 4.88163779  0.26349179  0.44771762 -0.14543905 -2.14512958 -0.03559434
 -0.29705992  4.26438321  0.74721543  0.38048825  2.09501546 -0.37111665
  1.81826542  0.74351787 -0.56926429 -0.59828605  0.51841252 -0.80646839]
RSS ERROR= 501.8642602243605
r=0.82, p=0.00

```

```

1 lr = LinearRegression()
2 lr.fit(newX,newY)
3 train_score, train_score_p = pearsonr(lr.predict(newX),newY)
4 print("\nr=" + format(train_score, '.2f') + ", p=" + format(train_score_p, '.2f'))
5
6 test_score, test_score_p = pearsonr(lr.predict(newX_test),newY_test)
7 print("r=" + format(test_score, '.2f') + ", p=" + format(test_score_p, '.2f'))

```

```

r=0.81, p=0.00
r=0.82, p=0.00

```

```

1 features_drop = ['Immoblized bus','Broken Truck', 'Running over', 'Occurrence\
2 involving freight', 'Incident involving dangerous freight', 'Fire','Manifestations']
3 betterX = train_data.drop(features_drop,axis=1)
4 betterX_test = test_data.drop(features_drop,axis=1)

```

```

1 #Creating an LR model using OLS Solution
2 dMtx = np.array(betterX)
3 dMtx = np.append(np.ones([100,1]),dMtx,axis=1)
4 dMtx = dMtx.astype('int64')
5
6 Outcome = np.array(newY)
7
8 first = (dMtx.T.dot(dMtx))
9 first = np.linalg.inv(first)
10 second = (dMtx.T.dot(newY))
11 w_opt = first.dot(second)
12 # print(w_opt) #this is the OLS model.
13
14 #Now using this model on our TEST DATA and calculate Pearson's r and RSS.
15

```

```

16 dMtx_test = np.array(betterX_test)
17 dMtx_test = np.append(np.ones([35,1]),dMtx_test,axis=1)
18 dMtx_test = dMtx_test.astype('int64')
19 test_outcome = np.array(newY_test)
20
21 RSS = 0
22 inner = dMtx_test.dot(w_opt)
23 first = (test_outcome-inner)
24 RSS = first.T.dot(first)
25
26 print('RSS ERROR=',RSS)
27 test_score, test_score_p = pearsonr(inner,test_outcome)
28 print("r=" + format(test_score, '.2f') + ", p=" + format(test_score_p, '.2f'))

```

```

❏ RSS ERROR= 3.9611861494767887e-26
   r=1.00, p=0.00

```

```

1 #using inbuilt functions just as to verify results.
2 print('\nUsing Scikit:')
3 lr = LinearRegression()
4 lr.fit(betterX,newY)
5
6 train_score, train_score_p = pearsonr(lr.predict(betterX),newY)
7 print("r=" + format(train_score, '.2f') + ", p=" + format(train_score_p, '.2f'))
8
9 test_score, test_score_p = pearsonr(lr.predict(betterX_test),newY_test)
10 print("r=" + format(test_score, '.2f') + ", p=" + format(test_score_p, '.2f'))

```

```

❏ Using Scikit:
   r=1.00, p=0.00
   r=1.00, p=0.00

```

```

1 values = []
2
3 for i in newY:
4     values.append(newY[i])
5 for i in newY_test:
6     values.append(newY_test[i])
7
8 summation = 0
9
10 for i in values:
11     summation += values[i]
12
13 mean = summation/len(values)
14
15 newNewY = []
16 for i in newY:
17     if (i>mean):
18         newNewY.append(1)
19     else:
20         newNewY.append(0)
21
22 newNewY_test = []
23 for i in newY_test:
24     if (i > mean):
25         newNewY_test.append(1)
26     else:
27         newNewY_test.append(0)
28

```

```

1 #logistic regression framework.
2 from sklearn.linear_model import LogisticRegression
3 from sklearn import preprocessing
4
5 newX_scaled = preprocessing.scale(newX)
6
7 LRModel = LogisticRegression(random_state=0).fit(newX_scaled,newNewY)
8 LRModel.predict(newX_test)
9 LRModel.score(newX_test,newNewY_test)
10

```

```

❏ 0.8857142857142857

```

```

1 #logistic regression framework.
2 betterX_scaled = preprocessing.scale(betterX)
3
4 LRModel = LogisticRegression(penalty='none').fit(betterX_scaled,newNewY)
5 LRModel.predict(betterX_test)
6 LRModel.score(betterX_test,newNewY_test)

```

```

❏ 0.8857142857142857

```

```

1 from sklearn.model_selection import LeaveOneOut
2 hyper = [0.0001, 0.001, 0.01, 0.1, 1, 10, 100, 1000, 10000]

```

```

3 loo = LeaveOneOut()
4 newX_scaled = np.array(newX_scaled)
5 newNewY = np.array(newNewY)
6 loo.get_n_splits(newX_scaled)
7 cross_val_accuracy = []
8 for h in hyper:
9     E = []
10    for train_index, test_index in loo.split(newX_scaled):
11        cross_train = []
12        cross_train = np.delete(newX_scaled, test_index, axis=0)
13        cross_train_outcome = np.delete(newNewY, test_index)
14        t = test_index.item
15        LRModel = LogisticRegression(C=h).fit(cross_train,cross_train_outcome)
16        LRModel.predict(newX_scaled[test_index])
17        E.append(LRModel.score(newX_scaled[test_index],newNewY[test_index]))
18    cross_val_accuracy.append(sum(E)/len(E))
19 print(cross_val_accuracy)
20 newX_test = preprocessing.scale(newX_test)
21 LRModel2 = LogisticRegression(C= 0.1).fit(newX_scaled,newNewY)
22 LRModel2.predict(newX_test)
23 LRModel2.score(newX_test,newNewY_test)

```

```

[0.57, 0.57, 0.61, 0.64, 0.6, 0.6, 0.59, 0.59, 0.59]
0.7428571428571429

```