

Lab-05_Unsupervised_Anomaly_Detection

321910302051

K.Kalyani

```
In [1]: import numpy as np
import pandas as pd
import seaborn as sns
```

```
In [2]: data=pd.read_csv("creditcard_lab_5.csv")
```

```
In [3]: data.shape
```

```
Out[3]: (284807, 31)
```

```
In [4]: data.head()
```

```
Out[4]:
```

	Time	V1	V2	V3	V4	V5	V6	V7	V8	V9	...	V21	V22	V23	V24	
0	0.0	-1.359807	-0.072781	2.536347	1.378155	-0.338321	0.462388	0.239599	0.098698	0.363787	...	-0.018307	0.277838	-0.110474	0.066928	0.
1	0.0	1.191857	0.266151	0.166480	0.448154	0.060018	-0.082361	-0.078803	0.085102	-0.255425	...	-0.225775	-0.638672	0.101288	-0.339846	0.
2	1.0	-1.358354	-1.340163	1.773209	0.379780	-0.503198	1.800499	0.791461	0.247676	-1.514654	...	0.247998	0.771679	0.909412	-0.689281	-0.
3	1.0	-0.966272	-0.185226	1.792993	-0.863291	-0.010309	1.247203	0.237609	0.377436	-1.387024	...	-0.108300	0.005274	-0.190321	-1.175575	0.
4	2.0	-1.158233	0.877737	1.548718	0.403034	-0.407193	0.095921	0.592941	-0.270533	0.817739	...	-0.009431	0.798278	-0.137458	0.141267	-0.

5 rows × 31 columns



```
In [5]: data.tail()
```

Out[5]:

	Time	V1	V2	V3	V4	V5	V6	V7	V8	V9	...	V21	V22	V23	
284802	172786.0	-11.881118	10.071785	-9.834783	-2.066656	-5.364473	-2.606837	-4.918215	7.305334	1.914428	...	0.213454	0.111864	1.014480	-0.50
284803	172787.0	-0.732789	-0.055080	2.035030	-0.738589	0.868229	1.058415	0.024330	0.294869	0.584800	...	0.214205	0.924384	0.012463	-1.01
284804	172788.0	1.919565	-0.301254	-3.249640	-0.557828	2.630515	3.031260	-0.296827	0.708417	0.432454	...	0.232045	0.578229	-0.037501	0.64
284805	172788.0	-0.240440	0.530483	0.702510	0.689799	-0.377961	0.623708	-0.686180	0.679145	0.392087	...	0.265245	0.800049	-0.163298	0.12
284806	172792.0	-0.533413	-0.189733	0.703337	-0.506271	-0.012546	-0.649617	1.577006	-0.414650	0.486180	...	0.261057	0.643078	0.376777	0.00

5 rows × 31 columns



In [6]: data.describe()

Out[6]:

	Time	V1	V2	V3	V4	V5	V6	V7	V8
count	284807.000000	2.848070e+05	2.848070e+05	2.848070e+05	2.848070e+05	2.848070e+05	2.848070e+05	2.848070e+05	2.848070e+05
mean	94813.859575	3.919560e-15	5.688174e-16	-8.769071e-15	2.782312e-15	-1.552563e-15	2.010663e-15	-1.694249e-15	-1.927028e-16
std	47488.145955	1.958696e+00	1.651309e+00	1.516255e+00	1.415869e+00	1.380247e+00	1.332271e+00	1.237094e+00	1.194353e+00
min	0.000000	-5.640751e+01	-7.271573e+01	-4.832559e+01	-5.683171e+00	-1.137433e+02	-2.616051e+01	-4.355724e+01	-7.321672e+01
25%	54201.500000	-9.203734e-01	-5.985499e-01	-8.903648e-01	-8.486401e-01	-6.915971e-01	-7.682956e-01	-5.540759e-01	-2.086297e-01
50%	84692.000000	1.810880e-02	6.548556e-02	1.798463e-01	-1.984653e-02	-5.433583e-02	-2.741871e-01	4.010308e-02	2.235804e-02
75%	139320.500000	1.315642e+00	8.037239e-01	1.027196e+00	7.433413e-01	6.119264e-01	3.985649e-01	5.704361e-01	3.273459e-01
max	172792.000000	2.454930e+00	2.205773e+01	9.382558e+00	1.687534e+01	3.480167e+01	7.330163e+01	1.205895e+02	2.000721e+01

8 rows × 31 columns



In [7]: data.info

```
Out[7]: <bound method DataFrame.info of
0      0.0 -1.359807 -0.072781  2.536347  1.378155 -0.338321
1      0.0  1.191857  0.266151  0.166480  0.448154  0.060018
2      1.0 -1.358354 -1.340163  1.773209  0.379780 -0.503198
```

```

3          1.0 -0.966272 -0.185226 1.792993 -0.863291 -0.010309
4          2.0 -1.158233 0.877737 1.548718 0.403034 -0.407193
...
284802 172786.0 -11.881118 10.071785 -9.834783 -2.066656 -5.364473
284803 172787.0 -0.732789 -0.055080 2.035030 -0.738589 0.868229
284804 172788.0 1.919565 -0.301254 -3.249640 -0.557828 2.630515
284805 172788.0 -0.240440 0.530483 0.702510 0.689799 -0.377961
284806 172792.0 -0.533413 -0.189733 0.703337 -0.506271 -0.012546

```

```

          V6      V7      V8      V9      ...      V21      V22  \
0      0.462388 0.239599 0.098698 0.363787 ... -0.018307 0.277838
1     -0.082361 -0.078803 0.085102 -0.255425 ... -0.225775 -0.638672
2      1.800499 0.791461 0.247676 -1.514654 ... 0.247998 0.771679
3      1.247203 0.237609 0.377436 -1.387024 ... -0.108300 0.005274
4      0.095921 0.592941 -0.270533 0.817739 ... -0.009431 0.798278
...
284802 -2.606837 -4.918215 7.305334 1.914428 ... 0.213454 0.111864
284803 1.058415 0.024330 0.294869 0.584800 ... 0.214205 0.924384
284804 3.031260 -0.296827 0.708417 0.432454 ... 0.232045 0.578229
284805 0.623708 -0.686180 0.679145 0.392087 ... 0.265245 0.800049
284806 -0.649617 1.577006 -0.414650 0.486180 ... 0.261057 0.643078

```

```

          V23      V24      V25      V26      V27      V28  Amount  \
0     -0.110474 0.066928 0.128539 -0.189115 0.133558 -0.021053 149.62
1      0.101288 -0.339846 0.167170 0.125895 -0.008983 0.014724 2.69
2      0.909412 -0.689281 -0.327642 -0.139097 -0.055353 -0.059752 378.66
3     -0.190321 -1.175575 0.647376 -0.221929 0.062723 0.061458 123.50
4     -0.137458 0.141267 -0.206010 0.502292 0.219422 0.215153 69.99
...
284802 1.014480 -0.509348 1.436807 0.250034 0.943651 0.823731 0.77
284803 0.012463 -1.016226 -0.606624 -0.395255 0.068472 -0.053527 24.79
284804 -0.037501 0.640134 0.265745 -0.087371 0.004455 -0.026561 67.88
284805 -0.163298 0.123205 -0.569159 0.546668 0.108821 0.104533 10.00
284806 0.376777 0.008797 -0.473649 -0.818267 -0.002415 0.013649 217.00

```

```

Class
0      0
1      0
2      0
3      0
4      0
...
284802 0
284803 0
284804 0
284805 0
284806 0

```

[284807 rows x 31 columns]>

```
In [8]: data.nunique()
```

```
Out[8]: Time      124592  
V1          275663  
V2          275663  
V3          275663  
V4          275663  
V5          275663  
V6          275663  
V7          275663  
V8          275663  
V9          275663  
V10         275663  
V11         275663  
V12         275663  
V13         275663  
V14         275663  
V15         275663  
V16         275663  
V17         275663  
V18         275663  
V19         275663  
V20         275663  
V21         275663  
V22         275663  
V23         275663  
V24         275663  
V25         275663  
V26         275663  
V27         275663  
V28         275663  
Amount      32767  
Class        2  
dtype: int64
```

```
In [9]: data['V10'].unique()
```

```
Out[9]: array([ 0.09079417, -0.16697441,  0.20764287, ..., -0.48478176,  
               -0.39912565, -0.91542665])
```

```
In [10]: data['Amount'].unique()
```

```
Out[10]: array([149.62,    2.69, 378.66, ..., 381.05, 337.54,  95.63])
```

```
In [11]: data.isnull().sum()
```

```
Out[11]: Time      0
         V1        0
         V2        0
         V3        0
         V4        0
         V5        0
         V6        0
         V7        0
         V8        0
         V9        0
         V10       0
         V11       0
         V12       0
         V13       0
         V14       0
         V15       0
         V16       0
         V17       0
         V18       0
         V19       0
         V20       0
         V21       0
         V22       0
         V23       0
         V24       0
         V25       0
         V26       0
         V27       0
         V28       0
         Amount    0
         Class     0
         dtype: int64
```

```
In [12]: data.isnull()
```

```
Out[12]:
```

	Time	V1	V2	V3	V4	V5	V6	V7	V8	V9	...	V21	V22	V23	V24	V25	V26	V27	V28	Amount	Class
0	False	False	False	False	False	False	False	False	False	False	...	False	False	False	False	False	False	False	False	False	False
1	False	False	False	False	False	False	False	False	False	False	...	False	False	False	False	False	False	False	False	False	False
2	False	False	False	False	False	False	False	False	False	False	...	False	False	False	False	False	False	False	False	False	False
3	False	False	False	False	False	False	False	False	False	False	...	False	False	False	False	False	False	False	False	False	False

	Time	V1	V2	V3	V4	V5	V6	V7	V8	V9	...	V21	V22	V23	V24	V25	V26	V27	V28	Amount	Class
4	False	False	False	False	False	False	False	False	False	False	...	False	False	False	False	False	False	False	False	False	False
...
284802	False	False	False	False	False	False	False	False	False	False	...	False	False	False	False	False	False	False	False	False	False
284803	False	False	False	False	False	False	False	False	False	False	...	False	False	False	False	False	False	False	False	False	False
284804	False	False	False	False	False	False	False	False	False	False	...	False	False	False	False	False	False	False	False	False	False
284805	False	False	False	False	False	False	False	False	False	False	...	False	False	False	False	False	False	False	False	False	False
284806	False	False	False	False	False	False	False	False	False	False	...	False	False	False	False	False	False	False	False	False	False

284807 rows × 31 columns

In [13]: `data.notnull().head()`

Out[13]:

	Time	V1	V2	V3	V4	V5	V6	V7	V8	V9	...	V21	V22	V23	V24	V25	V26	V27	V28	Amount	Class
0	True	True	True	True	True	True	True	True	True	True	...	True	True	True	True	True	True	True	True	True	True
1	True	True	True	True	True	True	True	True	True	True	...	True	True	True	True	True	True	True	True	True	True
2	True	True	True	True	True	True	True	True	True	True	...	True	True	True	True	True	True	True	True	True	True
3	True	True	True	True	True	True	True	True	True	True	...	True	True	True	True	True	True	True	True	True	True
4	True	True	True	True	True	True	True	True	True	True	...	True	True	True	True	True	True	True	True	True	True

5 rows × 31 columns

In [14]: `data.notnull().sum()`

Out[14]:

Time	284807
V1	284807
V2	284807
V3	284807
V4	284807
V5	284807
V6	284807
V7	284807

```

V8      284807
V9      284807
V10     284807
V11     284807
V12     284807
V13     284807
V14     284807
V15     284807
V16     284807
V17     284807
V18     284807
V19     284807
V20     284807
V21     284807
V22     284807
V23     284807
V24     284807
V25     284807
V26     284807
V27     284807
V28     284807
Amount  284807
Class   284807
dtype: int64

```

In [15]: `data.dropna()`

Out[15]:

	Time	V1	V2	V3	V4	V5	V6	V7	V8	V9	...	V21	V22	V23	
0	0.0	-1.359807	-0.072781	2.536347	1.378155	-0.338321	0.462388	0.239599	0.098698	0.363787	...	-0.018307	0.277838	-0.110474	C
1	0.0	1.191857	0.266151	0.166480	0.448154	0.060018	-0.082361	-0.078803	0.085102	-0.255425	...	-0.225775	-0.638672	0.101288	-C
2	1.0	-1.358354	-1.340163	1.773209	0.379780	-0.503198	1.800499	0.791461	0.247676	-1.514654	...	0.247998	0.771679	0.909412	-C
3	1.0	-0.966272	-0.185226	1.792993	-0.863291	-0.010309	1.247203	0.237609	0.377436	-1.387024	...	-0.108300	0.005274	-0.190321	-1
4	2.0	-1.158233	0.877737	1.548718	0.403034	-0.407193	0.095921	0.592941	-0.270533	0.817739	...	-0.009431	0.798278	-0.137458	C
...
284802	172786.0	-11.881118	10.071785	-9.834783	-2.066656	-5.364473	-2.606837	-4.918215	7.305334	1.914428	...	0.213454	0.111864	1.014480	-C
284803	172787.0	-0.732789	-0.055080	2.035030	-0.738589	0.868229	1.058415	0.024330	0.294869	0.584800	...	0.214205	0.924384	0.012463	-1
284804	172788.0	1.919565	-0.301254	-3.249640	-0.557828	2.630515	3.031260	-0.296827	0.708417	0.432454	...	0.232045	0.578229	-0.037501	C
284805	172788.0	-0.240440	0.530483	0.702510	0.689799	-0.377961	0.623708	-0.686180	0.679145	0.392087	...	0.265245	0.800049	-0.163298	C

	Time	V1	V2	V3	V4	V5	V6	V7	V8	V9	...	V21	V22	V23
284806	172792.0	-0.533413	-0.189733	0.703337	-0.506271	-0.012546	-0.649617	1.577006	-0.414650	0.486180	...	0.261057	0.643078	0.376777

284807 rows × 31 columns



In [16]: `data1=data.dropna()`

In [17]: `data1.isnull()`

Out[17]:

	Time	V1	V2	V3	V4	V5	V6	V7	V8	V9	...	V21	V22	V23	V24	V25	V26	V27	V28	Amount	Class
0	False	False	False	False	False	False	False	False	False	False	...	False	False	False	False	False	False	False	False	False	False
1	False	False	False	False	False	False	False	False	False	False	...	False	False	False	False	False	False	False	False	False	False
2	False	False	False	False	False	False	False	False	False	False	...	False	False	False	False	False	False	False	False	False	False
3	False	False	False	False	False	False	False	False	False	False	...	False	False	False	False	False	False	False	False	False	False
4	False	False	False	False	False	False	False	False	False	False	...	False	False	False	False	False	False	False	False	False	False
...
284802	False	False	False	False	False	False	False	False	False	False	...	False	False	False	False	False	False	False	False	False	False
284803	False	False	False	False	False	False	False	False	False	False	...	False	False	False	False	False	False	False	False	False	False
284804	False	False	False	False	False	False	False	False	False	False	...	False	False	False	False	False	False	False	False	False	False
284805	False	False	False	False	False	False	False	False	False	False	...	False	False	False	False	False	False	False	False	False	False
284806	False	False	False	False	False	False	False	False	False	False	...	False	False	False	False	False	False	False	False	False	False

284807 rows × 31 columns

In [18]: `import pandas as pd`
`import numpy as np`
`import math`
`import random`
`%matplotlib inline`
`import random`


```
from matplotlib import pyplot
import os
```

Anomaly Detection Algorithms: Isolation

Forest vs the Rest

Isolation Forests in Python

1. Forest
2. Isolation Tree
3. Evaluation (Path Length)

```
In [20]: class ExNode:
          def __init__(self, size):
              self.size = size

          class InNode:
              def __init__(self, left, right, splitAtt, splitVal):
                  self.left = left
                  self.right = right
                  self.splitAtt = splitAtt
                  self.splitVal = splitVal
```

FOREST

```
In [23]: def iForest(X, noOfTrees, sampleSize):
          forest = []
          hlim = math.ceil(math.log(sampleSize, 2))
          for i in range(noOfTrees):
              X_train = X.sample(sampleSize)
              forest.append(iTree(X_train, 0, hlim))
          return forest
```

ISOLATION TREE

```
In [24]: def iTree(X,currHeight,hlim):
          if currHeight>=hlim or len(X)<=1:
              return ExNode(len(X))
          else:
              Q=X.columns
              q=random.choice(Q)
              p=random.choice(X[q].unique())
              X_l=X[X[q]<p]
              X_r=X[X[q]>=p]
              return InNode(iTree(X_l,currHeight+1,hlim),iTree(X_r,currHeight+1,hlim),q,p)
```

PATH LENGTH

```
In [25]: def pathLength(x,Tree,currHeight):
          if isinstance(Tree,ExNode):
              return currHeight
          a=Tree.splitAtt
          if x[a]<Tree.splitVal:
              return pathLength(x,Tree.left,currHeight+1)
          else:
              return pathLength(x,Tree.right,currHeight+1)
```

TEST RUN

```
In [26]: df=pd.read_csv("creditcard_lab_5.csv")
          y_true=df['Class']
          df_data=df.drop('Class',1)
```

CREATING THE FOREST

```
In [27]: sampleSize=10000
          ifor=iForest(df_data.sample(100000),10,sampleSize)
```

Selecting 1000 random datapoints to get their path lengths. The purpose for this is to plot and see if anomalies actually have shorter path lengths.

```
In [28]: posLenLst=[]
negLenLst=[]
for sim in range(1000):
    ind=random.choice(df_data[y_true==1].index)
    for tree in ifor:
        posLenLst.append(pathLength(df_data.iloc[ind],tree,0))

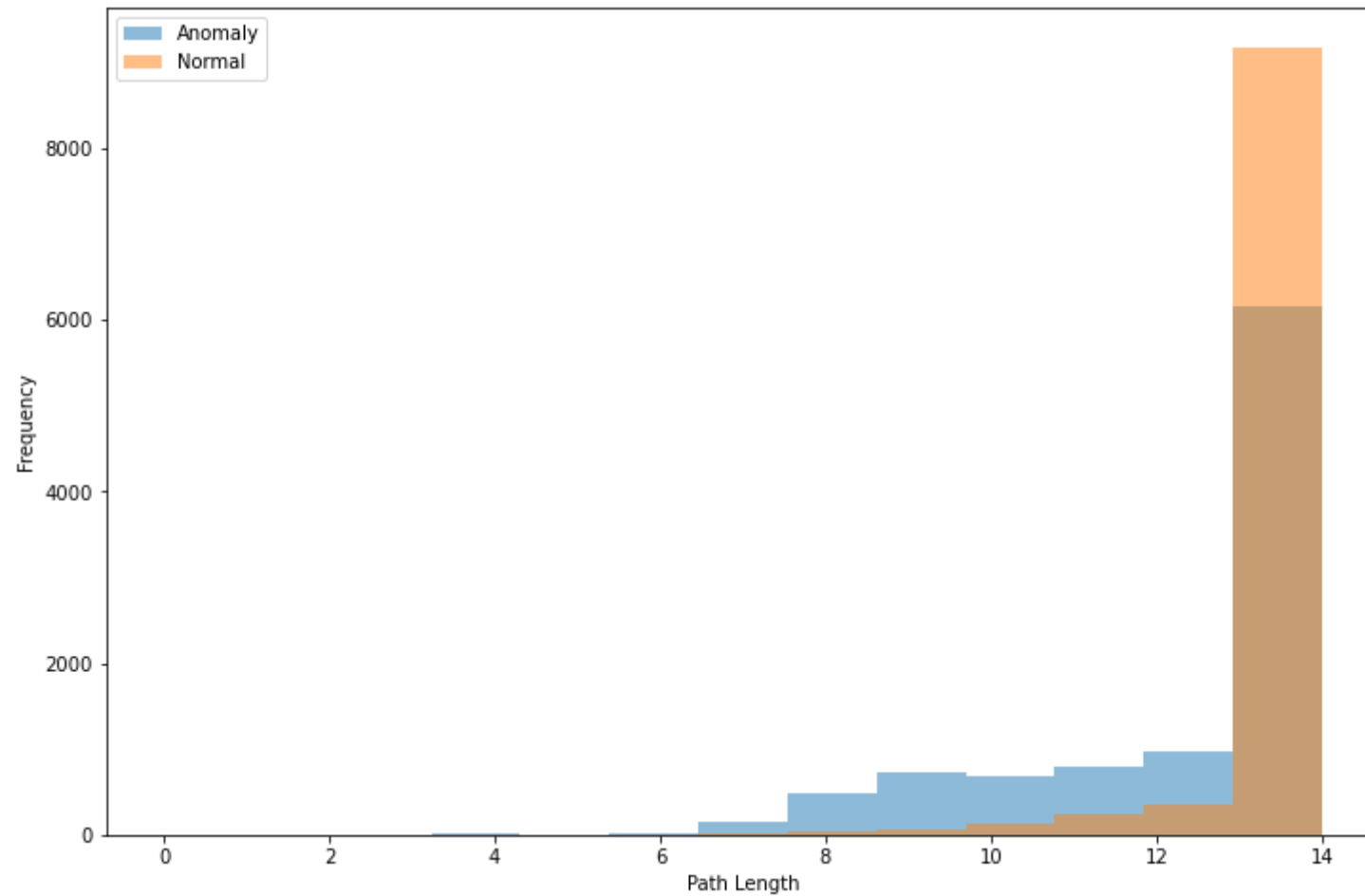
    ind=random.choice(df_data[y_true==0].index)
    for tree in ifor:
        negLenLst.append(pathLength(df_data.iloc[ind],tree,0))
```

Plotting the Path Lengths

```
In [32]: bins = np.linspace(0,math.ceil(math.log(sampleSize,2)), math.ceil(math.log(sampleSize,2)))

pyplot.figure(figsize=(12,8))
pyplot.hist(posLenLst, bins, alpha=0.5, label='Anomaly')
pyplot.hist(negLenLst, bins, alpha=0.5, label='Normal')
pyplot.xlabel('Path Length')
pyplot.ylabel('Frequency')
pyplot.legend(loc='upper left')
```

```
Out[32]: <matplotlib.legend.Legend at 0x1b92df25b20>
```



In []: