

```
In [ ]: CDSOFT
```

# Import Modules

```
In [1]: import pandas as pd
import numpy as np
import seaborn as sns
import matplotlib.pyplot as plt
import warnings
warnings.filterwarnings('ignore')
matplotlib inline
```

# Loading the Dataset

```
In [2]: df = pd.read_csv('creditcard.csv')
df.head()
```

	V1	V2	V3	V4	V5	V6	V7	V8	V9	...	V21	V22	V23	V24	V25	V26	V27	V28	Amount	
0	0.0	-1.359807	-0.072781	2.848070e+05	2.848070e+05	0.363787	...	...	...	-0.018307	0.277838	0.277838	0.277838	0.277838	0.277838	0.277838	0.277838	0.277838	149.6	
1	0.0	1.191857	0.266351	0.166480	0.448154	0.000018	-0.082361	-0.078803	0.085102	-0.255425	...	-0.225775	-0.388672	0.001288	-0.338646	0.167370	0.128899	-0.009883	0.034724	2.6
2	1.0	-1.359384	-1.340353	1.773209	0.379780	1.800499	0.791461	0.247676	-1.514654	...	0.247998	0.771679	0.909412	-0.699281	-0.327642	-0.198097	-0.055353	-0.059752	378.6	
3	1.0	-0.966272	-0.185226	1.782993	-0.863291	-0.003309	1.247203	0.237609	0.377436	-1.387024	...	-0.108300	0.005274	-0.190321	-1.175575	0.647376	-0.221829	0.062723	0.061458	123.9
4	2.0	-1.158233	0.877377	1.548718	0.403034	-0.407193	0.959221	0.929481	-0.270533	0.817739	...	-0.009431	0.798276	-0.137458	0.141267	-0.206010	0.502292	0.219422	0.215153	69.9

5 rows x 31 columns

```
In [3]: # statistical info
df.describe()
```

	Time	V1	V2	V3	V4	V5	V6	V7	V8	V9	...	V21	V22	
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	2.848070e+05	...	2.848070e+05	2.848070e+05	2.848070e+05
mean	14811.859575	1.168375e-15	3.416906e-16	-1.179537e-15	2.074095e-15	9.604066e-16	1.487313e-15	-5.556467e-16	1.213481e-16	0.066928	...	0.128939	-0.189115	0.133558
std	47481.145955	1.958566e+00	1.651309e+00	1.516255e+00	1.415899e+00	1.380247e+00	1.332271e+00	1.237094e+00	1.194533e+00	1.098632e+00	...	1.345240e+01	7.257016e+01	6.244601e+01
min	0.000000	-5.640751e-01	-7.371573e-01	-4.832559e-01	-5.663171e+00	-1.137433e+02	-2.616051e+01	-4.355724e+01	-7.321672e+01	-1.343407e+01	...	-3.483038e+01	-1.093314e+01	-4.480774e+01
25%	54201.500000	-9.203734e-01	-5.985499e-01	-8.903648e-01	-8.465401e-01	-6.915971e-01	-7.682956e-01	-5.540759e-01	-2.086297e-01	-6.430975e-02	...	-2.843949e+01	-5.423504e+01	-1.618461e+01
50%	84692.000000	1.110380e-02	6.548556e-02	1.798463e-01	1.984653e-02	5.433583e-02	2.741871e-01	4.010308e-02	2.238040e-02	-5.432973e-02	...	-2.945017e-02	6.781943e-03	-1.118293e-02
75%	139320.500000	1.315642e+00	8.837239e-01	1.027136e+00	7.433413e-01	6.119246e-01	3.985649e-01	5.704361e-01	3.273459e-01	5.971390e-01	...	1.863772e-01	5.285536e-01	1.476421e-01
max	172792.000000	2.443920e+00	2.220577e-01	9.382593e+00	1.687934e-01	3.480167e-01	7.330163e-01	1.205895e+02	2.000721e+01	1.559849e+01	...	2.720238e+01	1.050039e+01	2.252384e+01

8 rows x 31 columns

```
In [4]: # datatype info
df.info()
```

```
<class 'pandas.core.frame.DataFrame'>
RangeIndex: 284807 entries, 0 to 284806
Data columns (total 31 columns):
 #   Column  Non-Null Count  Dtype
---  --
 0   Time    284807 non-null      float64
 1   V1      284807 non-null      float64
 2   V2      284807 non-null      float64
 3   V3      284807 non-null      float64
 4   V4      284807 non-null      float64
 5   V5      284807 non-null      float64
 6   V6      284807 non-null      float64
 7   V7      284807 non-null      float64
 8   V8      284807 non-null      float64
 9   V9      284807 non-null      float64
10  V10     284807 non-null      float64
11  V11     284807 non-null      float64
12  V12     284807 non-null      float64
13  V13     284807 non-null      float64
14  V14     284807 non-null      float64
15  V15     284807 non-null      float64
16  V16     284807 non-null      float64
17  V17     284807 non-null      float64
18  V18     284807 non-null      float64
19  V19     284807 non-null      float64
20  V20     284807 non-null      float64
21  V21     284807 non-null      float64
22  V22     284807 non-null      float64
23  V23     284807 non-null      float64
24  V24     284807 non-null      float64
25  V25     284807 non-null      float64
26  V26     284807 non-null      float64
27  V27     284807 non-null      float64
28  V28     284807 non-null      float64
29  Amount  284807 non-null      float64
30  Class   284807 non-null      int64
31  ...     ...
memory usage: 67.4 MB
```

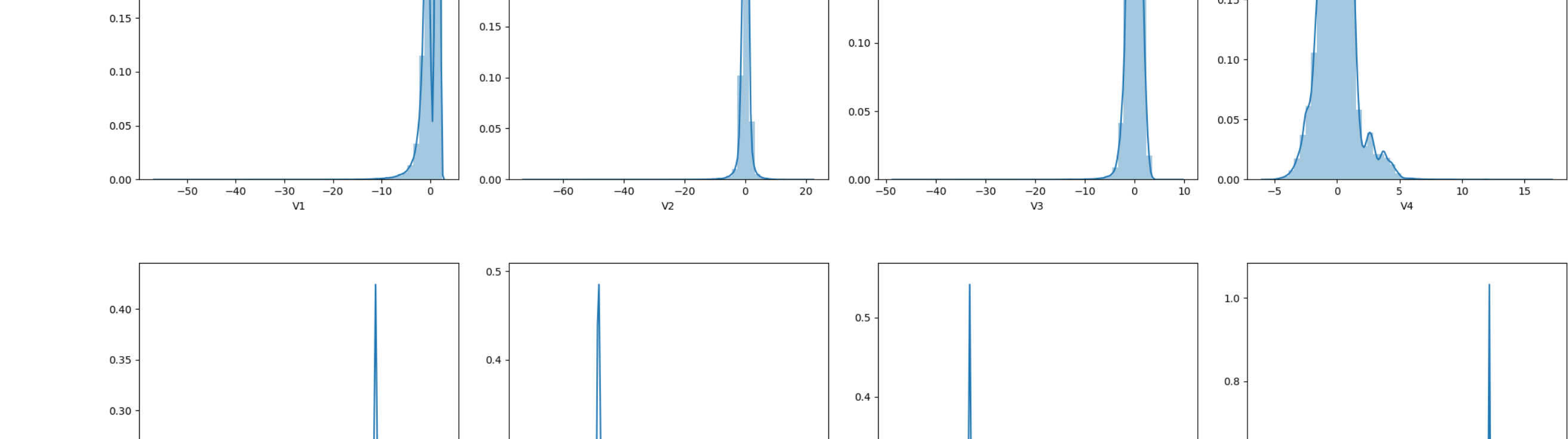
# Preprocessing the Dataset

```
In [5]: # check for null values
df.isnull().sum()
```

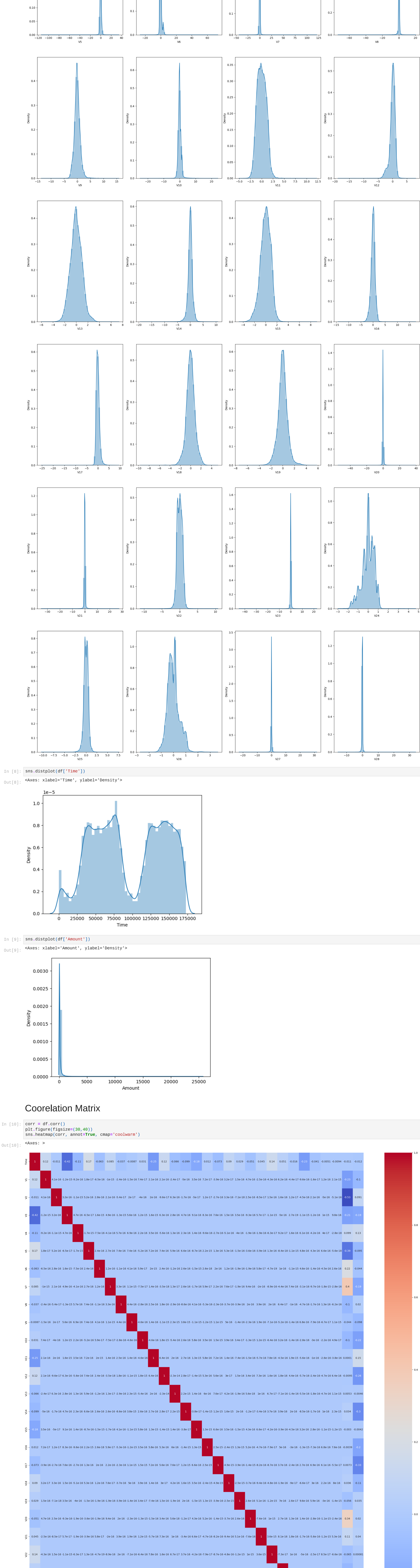
```
Out[5]: Time    0
V1         0
V2         0
V3         0
V4         0
V5         0
V6         0
V7         0
V8         0
V9         0
V10        0
V11        0
V12        0
V13        0
V14        0
V15        0
V16        0
V17        0
V18        0
V19        0
V20        0
V21        0
V22        0
V23        0
V24        0
V25        0
V26        0
V27        0
V28        0
Amount     0
Class      0
dtype: int64
```

# Exploratory Data Analysis

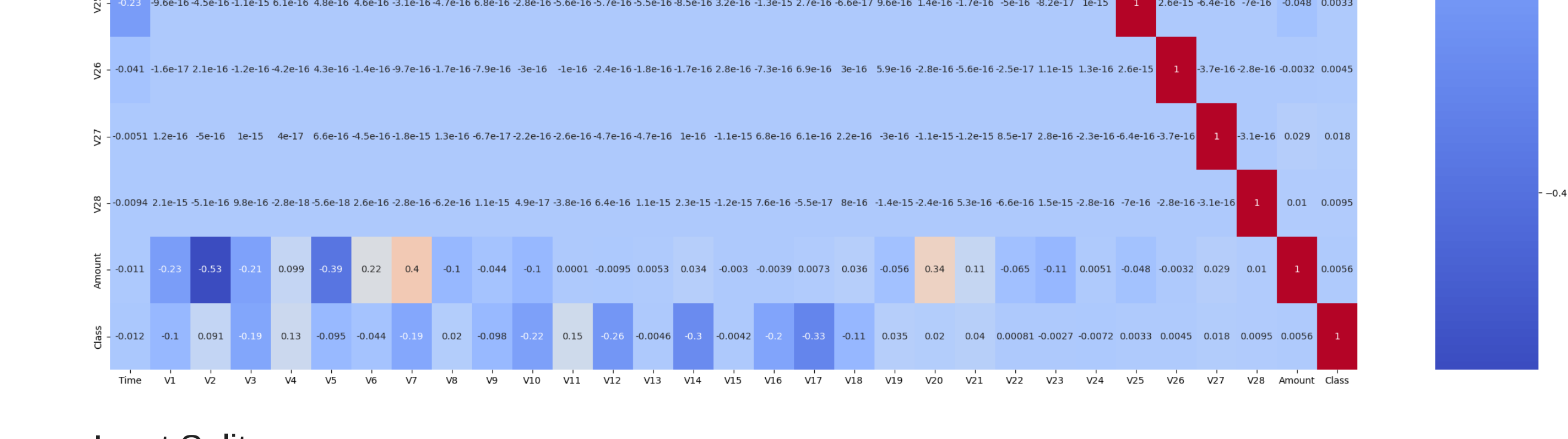
```
In [6]: sns.countplot(df['Class'])
<Axes: xlabel='count'>
```



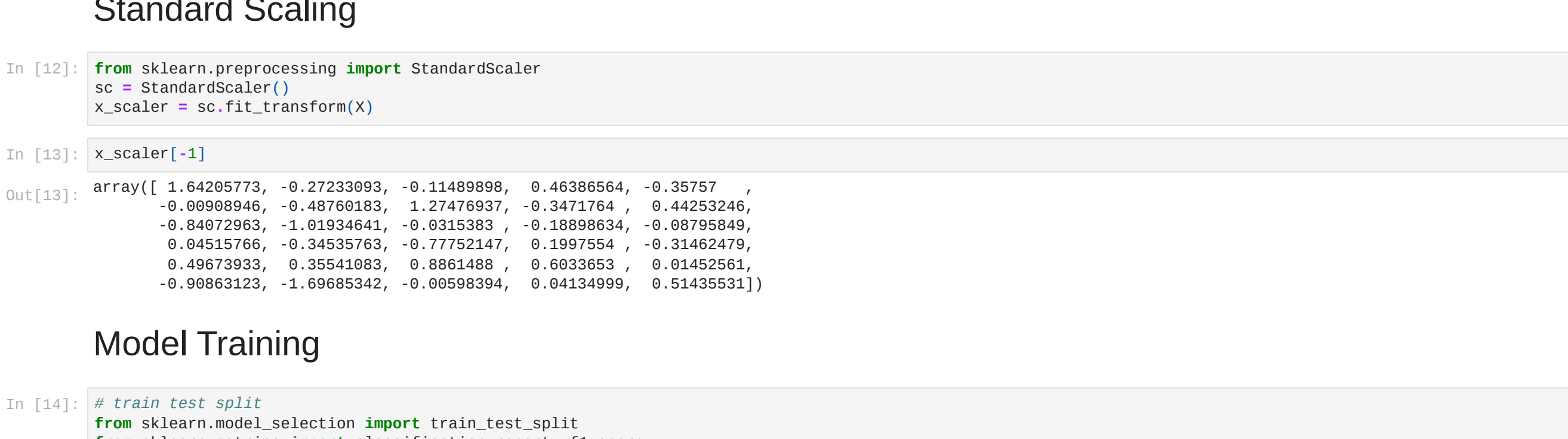
```
In [7]: df_temp = df.drop(columns=['Time', 'Amount', 'Class'], axis=1)
# create dist plots
fig, ax = plt.subplots(ncols=4, rows=7, figsize=(20, 50))
index = 0
ax = ax.flatten()
for col in df_temp.columns:
    sns.distplot(df_temp[col], ax=ax[index])
    index += 1
plt.tight_layout(pad=0.5, w_pad=0.5, h_pad=5)
```



```
In [8]: sns.distplot(df['Time'])
<Axes: xlabel='Time', ylabel='Density'>
```

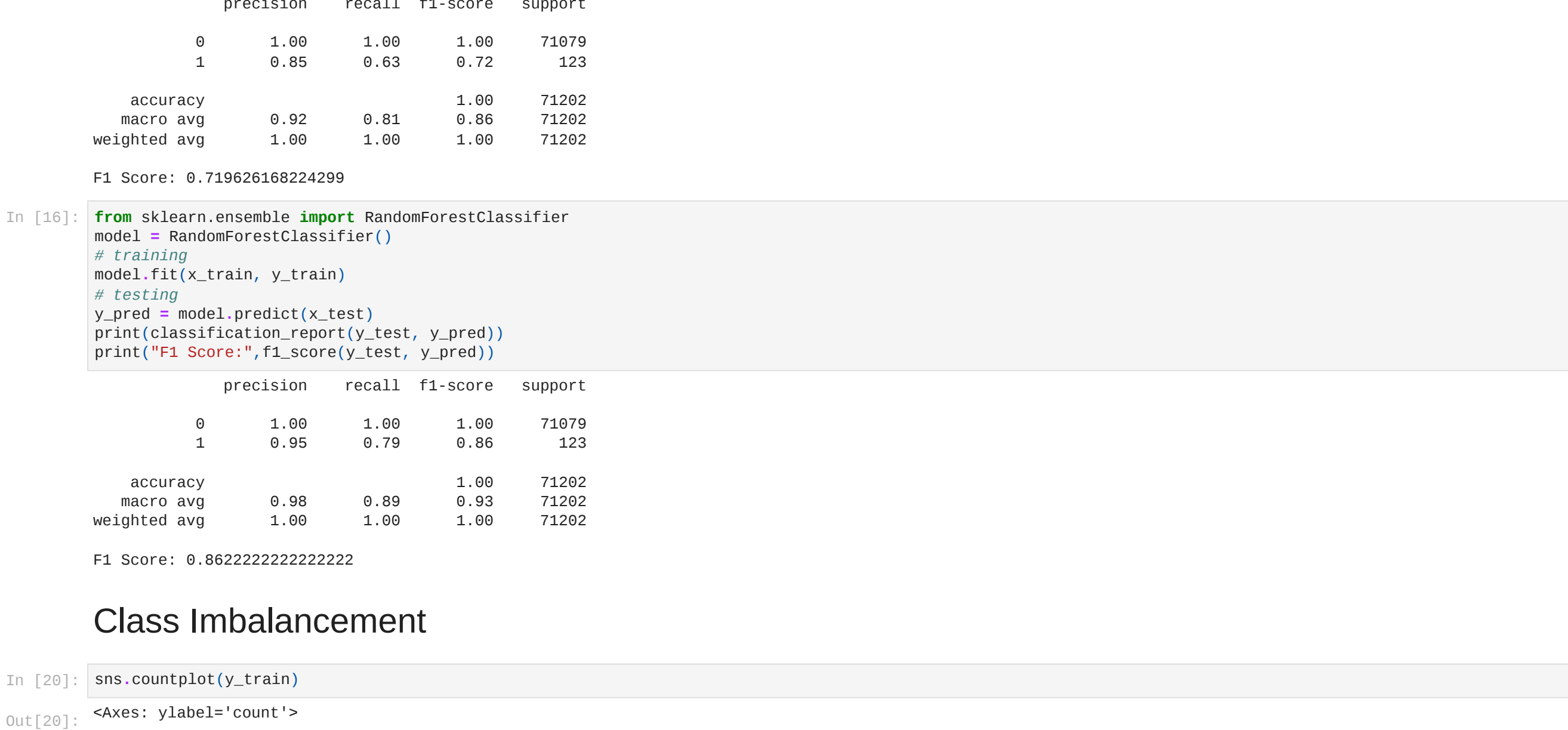


```
In [9]: sns.distplot(df['Amount'])
<Axes: xlabel='Amount', ylabel='Density'>
```



# Coorelation Matrix

```
In [10]: corr = df.corr()
plt.figure(figsize=(38,40))
sns.heatmap(corr, annot=True, cmap='coolwarm')
```



# Input Split

```
In [11]: X = df.drop(columns=['Class'], axis=1)
y = df['Class']
```

# Standard Scaling

```
In [12]: from sklearn.preprocessing import StandardScaler
sc = StandardScaler()
X_scaled = sc.fit_transform(X)
```

```
In [13]: X_scaled[1:]
```

```
Out[13]: array([[ 1.64265773, -0.27238993, -0.11489898,  0.46388664, -0.357575
-0.09898946, -0.48760353,  1.27478937, -0.3472764 ,  0.44535246,
-0.84072963, -1.01934481, -0.0315383 , -0.18888634, -0.08795849,
 0.64535766, -0.24832765, -0.7752147 ,  0.10979564, -0.31462479,
 0.49673933,  0.35541983,  0.8861488 ,  0.6033653 ,  0.01452561,
-0.00863123, -1.09685342, -0.09598394,  0.04134999,  0.51435531])
```

# Model Training

```
In [14]: # train test split
from sklearn.model_selection import train_test_split
from sklearn.metrics import classification_report, f1_score
x_train, x_test, y_train, y_test = train_test_split(X_scaled, y, test_size=0.25, random_state=42, stratify=y)
```

```
In [15]: from sklearn.linear_model import LogisticRegression
model = LogisticRegression()
# training
model.fit(x_train, y_train)
# testing
y_pred = model.predict(x_test)
print(classification_report(y_test, y_pred))
print("F1 Score:", f1_score(y_test, y_pred))
```

	precision	recall	f1-score	support
0	1.00	1.00	1.00	71079
1	0.85	0.63	0.72	123
accuracy	0.92	0.81	0.86	71202
macro avg	0.93	0.81	0.86	71202
weighted avg	1.00	0.99	1.00	71202
F1 Score:	0.71862618824299			

```
In [16]: from sklearn.ensemble import RandomForestClassifier
model = RandomForestClassifier()
# training
model.fit(x_train, y_train)
# testing
y_pred = model.predict(x_test)
print(classification_report(y_test, y_pred))
print("F1 Score:", f1_score(y_test, y_pred))
```

	precision	recall	f1-score	support
0	1.00	1.00	1.00	71079
1	0.95	0.80	0.83	123
accuracy	0.98	0.89	0.93	71202
macro avg	0.94	0.90	0.92	71202
weighted avg	1.00	0.99	0.99	71202
F1 Score:	0.834042553194993			

# Class Imbalancement

```
In [20]: sns.countplot(y_train)
<Axes: xlabel='count'>
```



```
In [28]: # hint - use combination of over sampling and under sampling
# balance the class with equal distribution
from imblearn.over_sampling import SMOTE
X_smote, y_smote = over_sample.fit_resample(x_train, y_train)
```

```
In [29]: sns.countplot(y_smote)
<Axes: xlabel='count'>
```



```
In [30]: from sklearn.linear_model import LogisticRegression
model = LogisticRegression()
# training
model.fit(X_smote, y_smote)
# testing
y_pred = model.predict(x_test)
print(classification_report(y_test, y_pred))
print("F1 Score:", f1_score(y_test, y_pred))
```

	precision	recall	f1-score	support
0	1.00	1.00	1.00	71079
1	0.85	0.89	0.87	123
accuracy	0.92	0.93	0.92	71202
macro avg	0.94	0.96	0.95	71202
weighted avg	1.00	0.99	0.99	71202
F1 Score:	0.11621438812084			

```
In [31]: from sklearn.ensemble import RandomForestClassifier
model = RandomForestClassifier(n_jobs=-1)
# training
model.fit(X_smote, y_smote)
# testing
y_pred = model.predict(x_test)
print(classification_report(y_test, y_pred))
print("F1 Score:", f1_score(y_test, y_pred))
```

	precision	recall	f1-score	support
0	1.00	1.00	1.00	71079
1	0.88	0.80	0.83	123
accuracy	0.98	0.89	0.93	71202
macro avg	0.94	0.90	0.92	71202
weighted avg	1.00	0.99	0.99	71202
F1 Score:	0.834042553194993			