

```
In [2]: import pandas as pd
import numpy as np
from mpl_toolkits.mplot3d import Axes3D
from sklearn.preprocessing import StandardScaler
import matplotlib.pyplot as plt
import os
```

```
In [3]: df = pd.read_csv('creditcard.csv')
df.head()
```

Out[3]:

	Time	V1	V2	V3	V4	V5	V6	V7	V8	V9	...	
0	0.0	-1.359807	-0.072781	2.536347	1.378155	-0.338321	0.462388	0.239599	0.098698	0.363787	...	-0.
1	0.0	1.191857	0.266151	0.166480	0.448154	0.060018	-0.082361	-0.078803	0.085102	-0.255425	...	-0.
2	1.0	-1.358354	-1.340163	1.773209	0.379780	-0.503198	1.800499	0.791461	0.247676	-1.514654	...	0.
3	1.0	-0.966272	-0.185226	1.792993	-0.863291	-0.010309	1.247203	0.237609	0.377436	-1.387024	...	-0.
4	2.0	-1.158233	0.877737	1.548718	0.403034	-0.407193	0.095921	0.592941	-0.270533	0.817739	...	-0.

5 rows × 31 columns

```
In [4]: df.shape
```

Out[4]: (284807, 31)

```
In [5]: df.describe()
```

Out[5]:

	Time	V1	V2	V3	V4	V5	V
count	284807.000000	2.848070e+05	2.848070e+05	2.848070e+05	2.848070e+05	2.848070e+05	2.848070e+0
mean	94813.859575	1.168375e-15	3.416908e-16	-1.379537e-15	2.074095e-15	9.604066e-16	1.487313e-1
std	47488.145955	1.958696e+00	1.651309e+00	1.516255e+00	1.415869e+00	1.380247e+00	1.332271e+0
min	0.000000	-5.640751e+01	-7.271573e+01	-4.832559e+01	-5.683171e+00	-1.137433e+02	-2.616051e+0
25%	54201.500000	-9.203734e-01	-5.985499e-01	-8.903648e-01	-8.486401e-01	-6.915971e-01	-7.682956e-0
50%	84692.000000	1.810880e-02	6.548556e-02	1.798463e-01	-1.984653e-02	-5.433583e-02	-2.741871e-0
75%	139320.500000	1.315642e+00	8.037239e-01	1.027196e+00	7.433413e-01	6.119264e-01	3.985649e-0
max	172792.000000	2.454930e+00	2.205773e+01	9.382558e+00	1.687534e+01	3.480167e+01	7.330163e+0

8 rows × 31 columns

```
In [6]: 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
```

```

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
dtypes: float64(30), int64(1)
memory usage: 67.4 MB

```

In [7]: `df.isnull().sum()`

```

Out[7]: 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

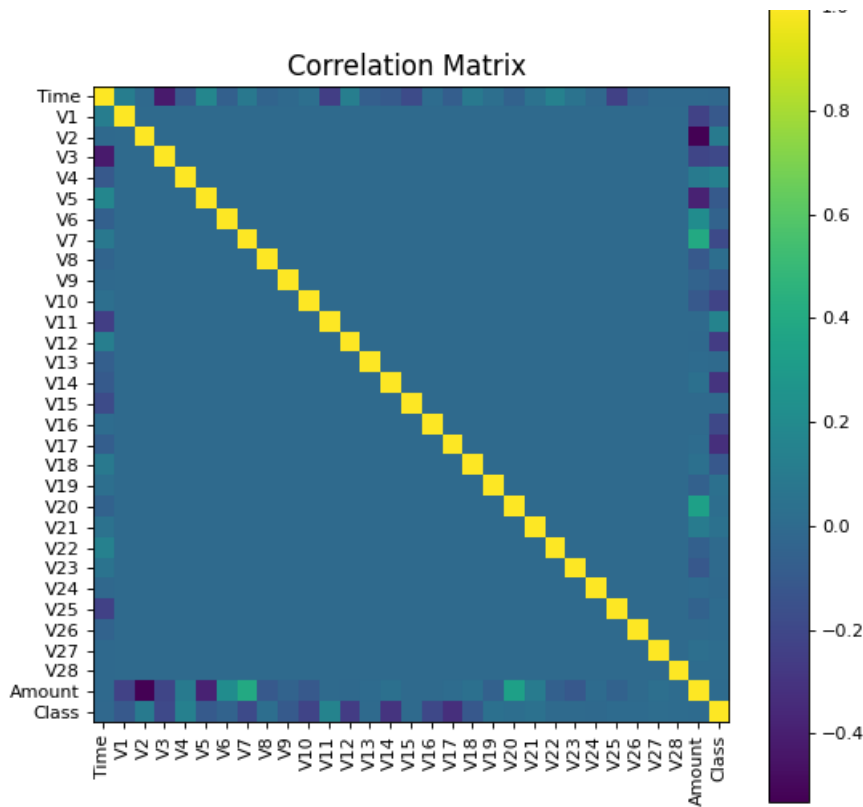
```

```

In [8]: # Correlation matrix
def plotCorrelationMatrix(df, graphWidth):
    # filename = df.dataframeName
    # df = df.dropna('columns') # drop columns with NaN
    df = df[[col for col in df if df[col].nunique() > 1]] # keep columns where there are more than 1
    if df.shape[1] < 2:
        print(f'No correlation plots shown: The number of non-NaN or constant columns ({df.shape[1]})
        return
    corr = df.corr()
    plt.figure(num=None, figsize=(graphWidth, graphWidth), dpi=80, facecolor='w', edgecolor='k')
    corrMat = plt.matshow(corr, fignum = 1)
    plt.xticks(range(len(corr.columns)), corr.columns, rotation=90)
    plt.yticks(range(len(corr.columns)), corr.columns)
    plt.gca().xaxis.tick_bottom()
    plt.colorbar(corrMat)
    plt.title(f'Correlation Matrix ', fontsize=15)
    plt.show()

```

In [9]: `plotCorrelationMatrix(df, 8)`



```
In [10]: # Count how many times each data type is present in the dataset
pd.value_counts(df.dtypes)
```

```
Out[10]: float64    30
         int64      1
         Name: count, dtype: int64
```

```
In [11]: df.isna()
```

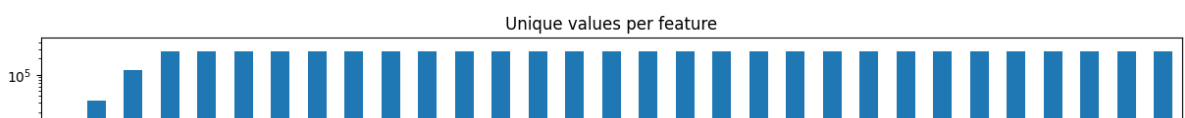
```
Out[11]:
```

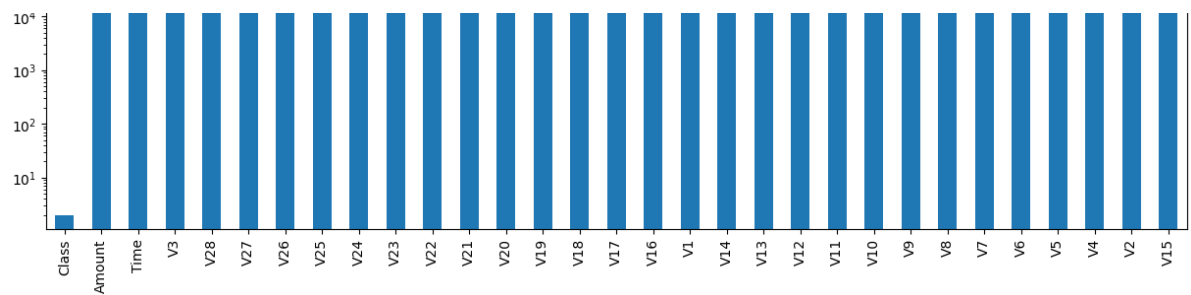
	Time	V1	V2	V3	V4	V5	V6	V7	V8	V9	...	V21	V22	V23	V24	V25	V26
0	False	False	False	False	False	False	False	False	False	False	...	False	False	False	False	False	False
1	False	False	False	False	False	False	False	False	False	False	...	False	False	False	False	False	False
2	False	False	False	False	False	False	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
4	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
284803	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
284805	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

284807 rows × 31 columns

```
In [12]: unique_values = df.select_dtypes(include="number").nunique().sort_values()
         unique_values.plot(logy=True, figsize=(15, 4), title="Unique values per feature")
```

```
Out[12]: <Axes: title={'center': 'Unique values per feature'}>
```

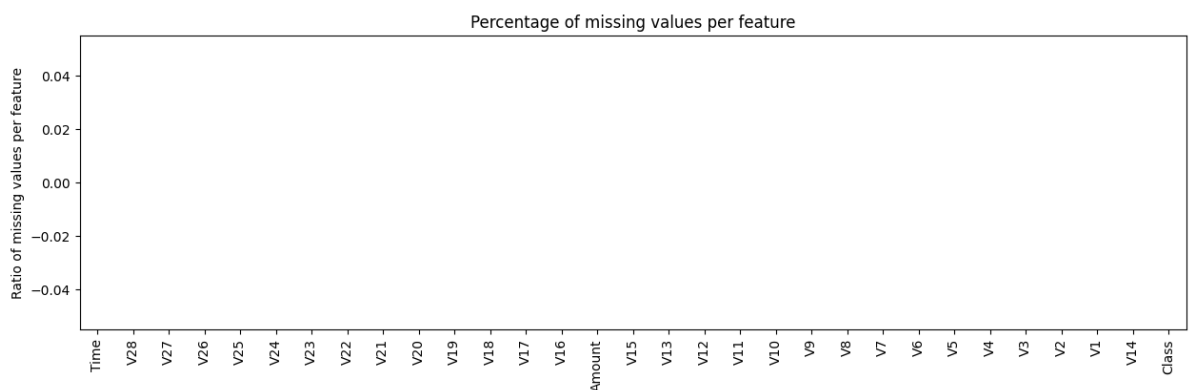




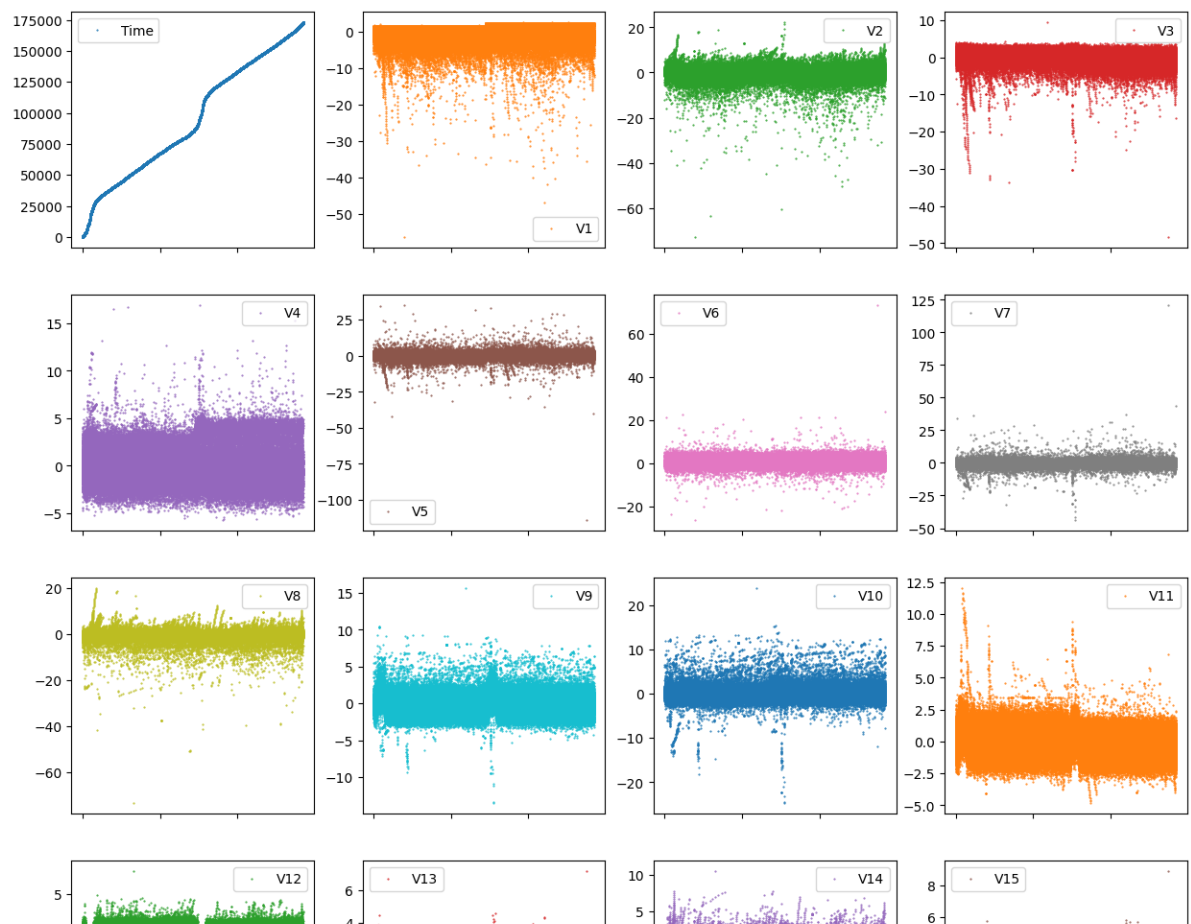
```
In [15]: n_duplicates = df.drop(labels=["Class"], axis=1).duplicated().sum()
print(f"You seem to have {n_duplicates} duplicates in your database.")
```

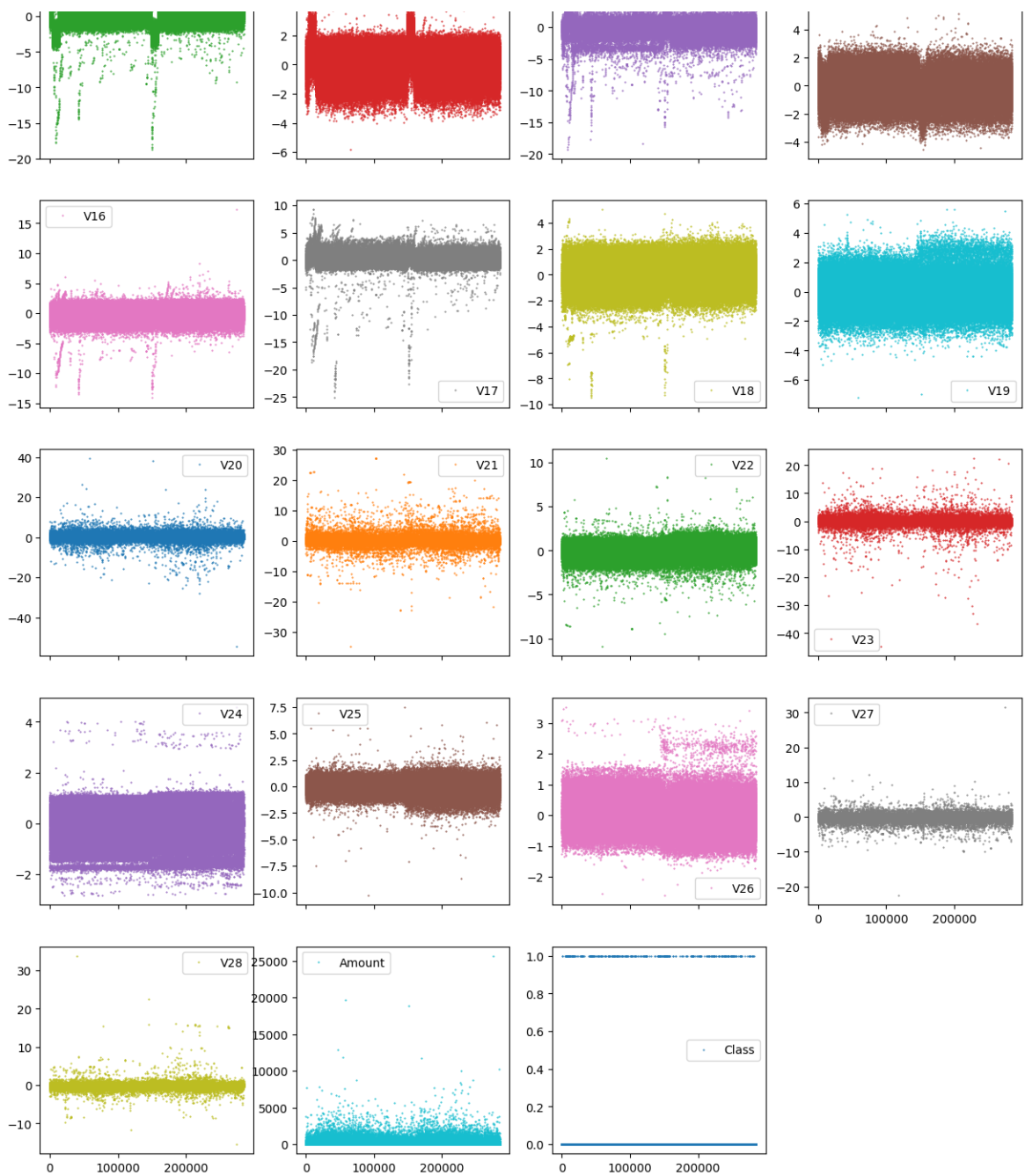
You seem to have 1081 duplicates in your database.

```
In [16]: df.isna().mean().sort_values().plot(
    kind="bar", figsize=(15, 4),
    title="Percentage of missing values per feature",
    ylabel="Ratio of missing values per feature");
#that is no missing value
```

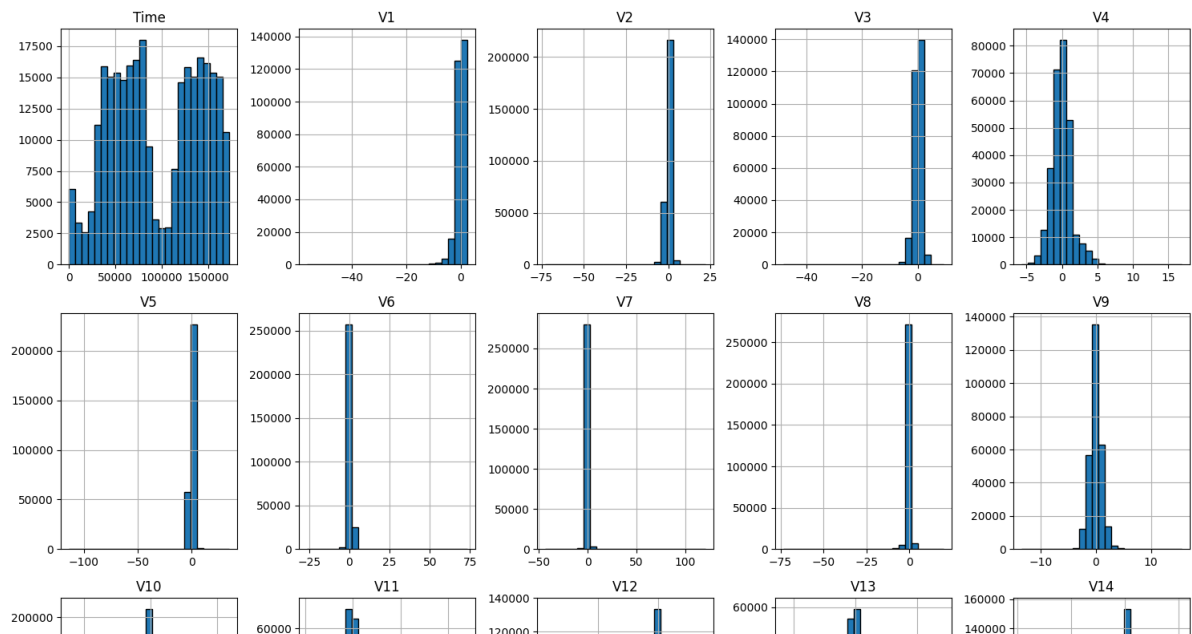


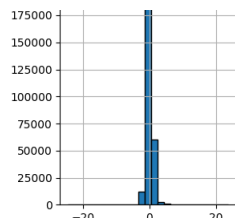
```
In [17]: df.plot(lw=0, marker=".", subplots=True, layout=(-1, 4),
    figsize=(15, 30), markersize=1);
```



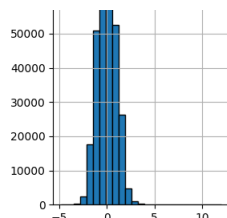


```
In [21]: df.hist(bins=25, figsize=(15, 25), layout=(-1, 5), edgecolor="black")
plt.tight_layout();
```

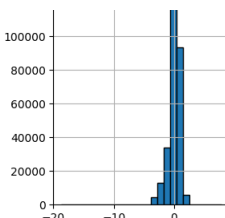




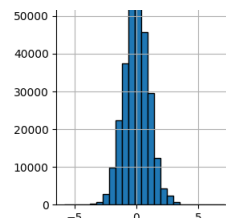
V15



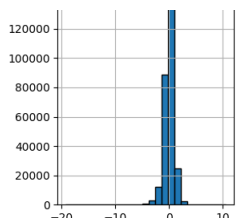
V16



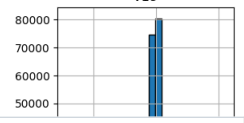
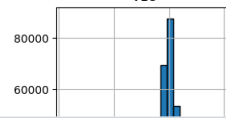
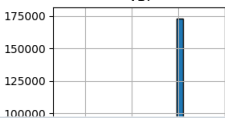
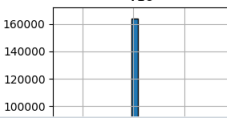
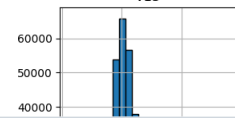
V17



V18



V19



```
In [1]: import numpy as np # linear algebra
import pandas as pd # data processing, CSV file I/O (e.g. pd.read_csv)
import tensorflow as tf
import matplotlib.pyplot as plt
import seaborn as sns
from sklearn.manifold import TSNE
from sklearn.decomposition import PCA, TruncatedSVD
import matplotlib.patches as mpatches
import time
```

```
# Classifier Libraries
from sklearn.linear_model import LogisticRegression
from sklearn.svm import SVC
from sklearn.neighbors import KNeighborsClassifier
from sklearn.tree import DecisionTreeClassifier
from sklearn.ensemble import RandomForestClassifier
import collections
```

```
In [ ]: df = pd.read_csv("creditcard.csv")
```

```
In [ ]: Total_transactions = len(df)
normal = len(df[df.Class == 0])
fraudulent = len(df[df.Class == 1])
fraud_percentage = round(fraudulent/normal*100, 2)

print("total number of transactions: ",Total_transactions)
print("normal transactions: ",normal)
print("fraud transactions: ",fraudulent)
print("percentage: ",fraud_percentage)
```

```
total number of transactions: 284807
normal transactions: 284315
fraud transactions: 492
percentage: 0.17
```

```
In [ ]: df.head()
```

```
Out[ ]:
```

	Time	V1	V2	V3	V4	V5	V6	V7	V8	V9	...
0	0.0	-1.359807	-0.072781	2.536347	1.378155	-0.338321	0.462388	0.239599	0.098698	0.363787	...
1	0.0	1.191857	0.266151	0.166480	0.448154	0.060018	-0.082361	-0.078803	0.085102	-0.255425	...
2	1.0	-1.358354	-1.340163	1.773209	0.379780	-0.503198	1.800499	0.791461	0.247676	-1.514654	...
3	1.0	-0.966272	-0.185226	1.792993	-0.863291	-0.010309	1.247203	0.237609	0.377436	-1.387024	...
4	2.0	-1.158233	0.877737	1.548718	0.403034	-0.407193	0.095921	0.592941	-0.270533	0.817739	...

5 rows × 31 columns

```
In [ ]: df.isnull().sum()
```

```
Out[ ]: 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 [ ]: print(df.duplicated().sum())
        df.shape
```

1081

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

```
In [ ]: df.drop_duplicates()
```

```
Out[ ]:
```

	scaled_amount	scaled_time	V1	V2	V3	V4	V5	V6	V7
110192	0.162230	-0.152434	-0.479760	0.902012	1.678826	0.679300	-0.383552	-0.989120	0.716873
82178	-0.027947	-0.298570	-1.100552	-0.646598	2.384530	-1.618709	-0.307001	0.379083	-0.865919
190623	-0.172990	0.519496	-0.986606	1.030981	0.635763	-1.308280	1.520850	-0.202459	1.226071
48178	0.076853	-0.483617	-1.099733	-0.031313	1.985752	0.380360	-0.885633	0.166664	-0.957074
240411	14.205268	0.774222	0.169344	-3.931869	-3.761457	-0.747751	-0.342344	-0.151290	1.720978
...
85825	-0.293440	-0.278868	-0.989250	0.651069	-0.998413	-1.408951	2.812663	3.547332	-0.818144
169348	0.321246	0.410719	-0.596156	-0.105169	0.846754	-1.464162	0.637215	-0.567283	0.527580
104993	-0.167819	-0.180594	1.155033	-0.192103	0.433161	0.662576	-0.333750	0.358379	-0.389224
66337	-0.294977	-0.384074	0.166544	1.131385	0.144757	0.743303	0.594937	-0.621800	0.917669
115546	-0.234053	-0.126728	-0.801655	1.378684	0.641833	-0.324040	0.691111	0.010191	0.590762

283726 rows × 31 columns

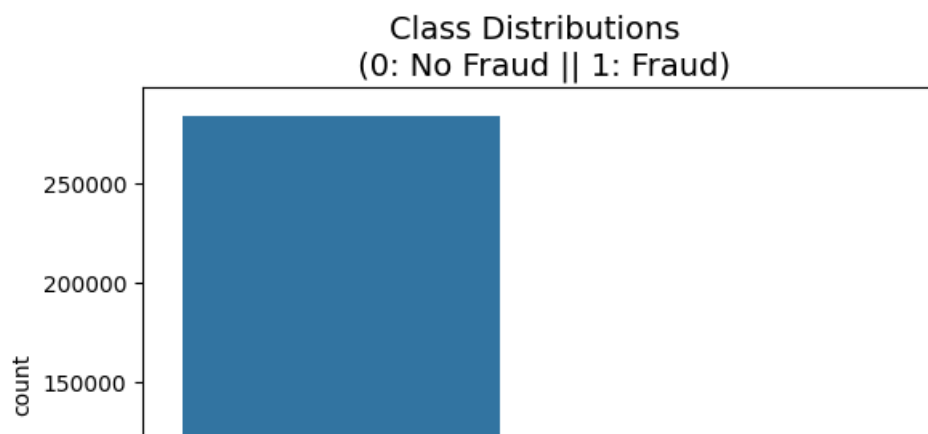


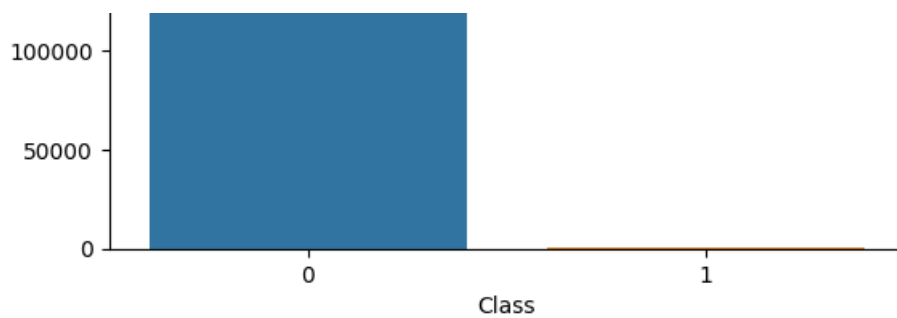
Pre-Processing

```
In [ ]: import matplotlib.pyplot as plt
        colors = ["#0101DF", "#DF0101"]

        sns.countplot(x='Class', data=df)
        plt.title('Class Distributions \n (0: No Fraud || 1: Fraud)', fontsize=14)
```

```
Out[ ]: Text(0.5, 1.0, 'Class Distributions \n (0: No Fraud || 1: Fraud)')
```





```
In [ ]: import warnings
warnings.filterwarnings("ignore")
```

```
In [ ]: fig, ax = plt.subplots(1, 2, figsize=(18,4))

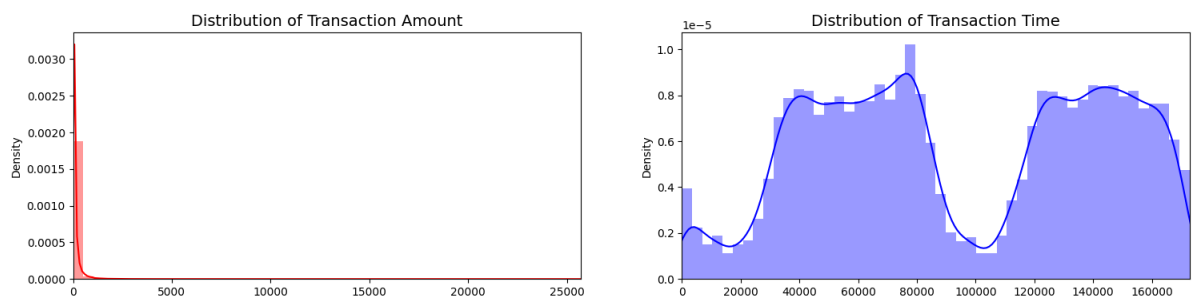
amount_val = df['Amount'].values
time_val = df['Time'].values

sns.distplot(amount_val, ax=ax[0], color='r')
ax[0].set_title('Distribution of Transaction Amount', fontsize=14)
ax[0].set_xlim([min(amount_val), max(amount_val)])

sns.distplot(time_val, ax=ax[1], color='b')
ax[1].set_title('Distribution of Transaction Time', fontsize=14)
ax[1].set_xlim([min(time_val), max(time_val)])

plt.show()

#very imbalanced data in amount and time
```



```
In [ ]: #scaling columns amount and time

from sklearn.preprocessing import StandardScaler, RobustScaler

# RobustScaler

std_scaler = StandardScaler()
rob_scaler = RobustScaler()

df['scaled_amount'] = rob_scaler.fit_transform(df['Amount'].values.reshape(-1,1))
df['scaled_time'] = rob_scaler.fit_transform(df['Time'].values.reshape(-1,1))

df.drop(['Time', 'Amount'], axis=1, inplace=True)
```

```
In [ ]: scaled_amount = df['scaled_amount']
scaled_time = df['scaled_time']

df.drop(['scaled_amount', 'scaled_time'], axis=1, inplace=True)
df.insert(0, 'scaled_amount', scaled_amount)
df.insert(1, 'scaled_time', scaled_time)

# Amount and Time are Scaled!

df.head()
```

```
Out [ ]: scaled_amount  scaled_time    V1    V2    V3    V4    V5    V6    V7
0          1.783274   -0.994983 -1.359807 -0.072781  2.536347  1.378155 -0.338321  0.462388  0.239599  0.098
```

1	-0.269825	-0.994983	1.191857	0.266151	0.166480	0.448154	0.060018	-0.082361	-0.078803	0.085
2	4.983721	-0.994972	-1.358354	-1.340163	1.773209	0.379780	-0.503198	1.800499	0.791461	0.247
3	1.418291	-0.994972	-0.966272	-0.185226	1.792993	-0.863291	-0.010309	1.247203	0.237609	0.377
4	0.670579	-0.994960	-1.158233	0.877737	1.548718	0.403034	-0.407193	0.095921	0.592941	-0.270

5 rows × 31 columns

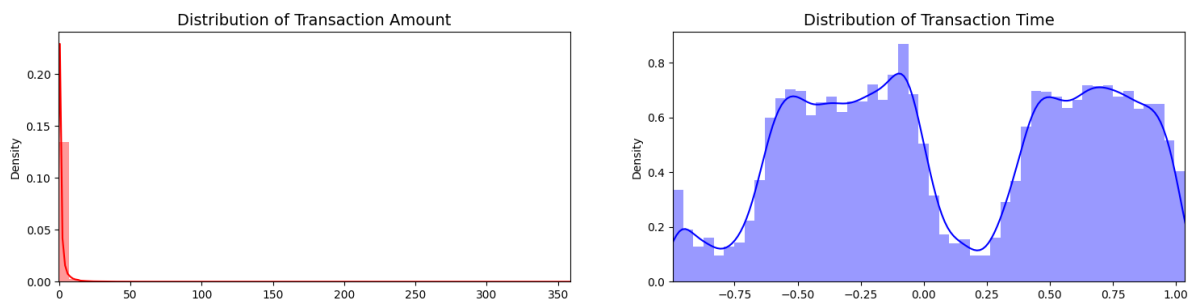
```
In [ ]: fig, ax = plt.subplots(1, 2, figsize=(18,4))

amount_val = df['scaled_amount'].values
time_val = df['scaled_time'].values

sns.distplot(amount_val, ax=ax[0], color='r')
ax[0].set_title('Distribution of Transaction Amount', fontsize=14)
ax[0].set_xlim([min(amount_val), max(amount_val)])

sns.distplot(time_val, ax=ax[1], color='b')
ax[1].set_title('Distribution of Transaction Time', fontsize=14)
ax[1].set_xlim([min(time_val), max(time_val)])

plt.show()
```



```
In [ ]: from sklearn.model_selection import train_test_split
from sklearn.model_selection import StratifiedShuffleSplit
from sklearn.model_selection import KFold, StratifiedKFold

print('No Frauds', round(df['Class'].value_counts()[0]/len(df) * 100,2), '% of the dataset')
print('Frauds', round(df['Class'].value_counts()[1]/len(df) * 100,2), '% of the dataset')

X = df.drop('Class', axis=1)
y = df['Class']

sss = StratifiedKFold(n_splits=5, random_state=None, shuffle=False)

for train_index, test_index in sss.split(X, y):
    print("Train:", train_index, "Test:", test_index)
    original_Xtrain, original_Xtest = X.iloc[train_index], X.iloc[test_index]
    original_ytrain, original_ytest = y.iloc[train_index], y.iloc[test_index]
```

```
No Frauds 99.83 % of the dataset
Frauds 0.17 % of the dataset
Train: [ 30473  30496  31002 ... 284804 284805 284806] Test: [    0     1     2 ... 57017 57018 57019]
Train: [    0     1     2 ... 284804 284805 284806] Test: [ 30473  30496  31002 ... 113964 113965 113966]
Train: [    0     1     2 ... 284804 284805 284806] Test: [ 81609  82400  83053 ... 170946 170947 170948]
Train: [    0     1     2 ... 284804 284805 284806] Test: [150654 150660 150661 ... 227866 227867 227868]
Train: [    0     1     2 ... 227866 227867 227868] Test: [212516 212644 213092 ... 284804 284805 284806]
```

```
In [ ]: # Random Under Sampling" which basically consists of removing data in order to have a more balanced d
# Lets shuffle the data before creating the subsamples

df = df.sample(frac=1)

# amount of fraud classes 492 rows.
fraud_df = df.loc[df['Class'] == 1]
```

```

fraud_df = df.loc[df['Class'] == 1]
non_fraud_df = df.loc[df['Class'] == 0][:492]

normal_distributed_df = pd.concat([fraud_df, non_fraud_df])

# Shuffle dataframe rows
new_df = normal_distributed_df.sample(frac=1, random_state=42)

new_df.head()

```

```

Out [ ]:

```

	scaled_amount	scaled_time	V1	V2	V3	V4	V5	V6	V
151906	-0.046112	0.138359	2.057176	0.180285	-1.870319	0.187160	0.523782	-1.505314	0.66731
249963	-0.296653	0.821967	-0.679521	4.672553	-6.814798	7.143500	0.928654	-1.873013	-2.30668
110046	0.314120	-0.153115	1.209445	0.085729	-0.048025	-0.093705	-0.082630	-0.772552	0.31248
150654	-0.307273	0.107403	-3.765680	5.890735	-10.202268	10.259036	-5.611448	-3.235376	-10.63268
154454	1.758821	0.198604	0.913116	1.145381	-4.602878	2.091803	-0.473224	-2.085436	-1.67124

5 rows × 31 columns

```

In [ ]:
print('Distribution of the Classes in the subsample dataset')
print(new_df['Class'].value_counts()/len(new_df))

sns.countplot(x='Class', data=new_df, palette=colors)
plt.title('Equally Distributed Classes', fontsize=14)
plt.show()

```

Distribution of the Classes in the subsample dataset

Class

0 0.5

1 0.5

Name: count, dtype: float64



```

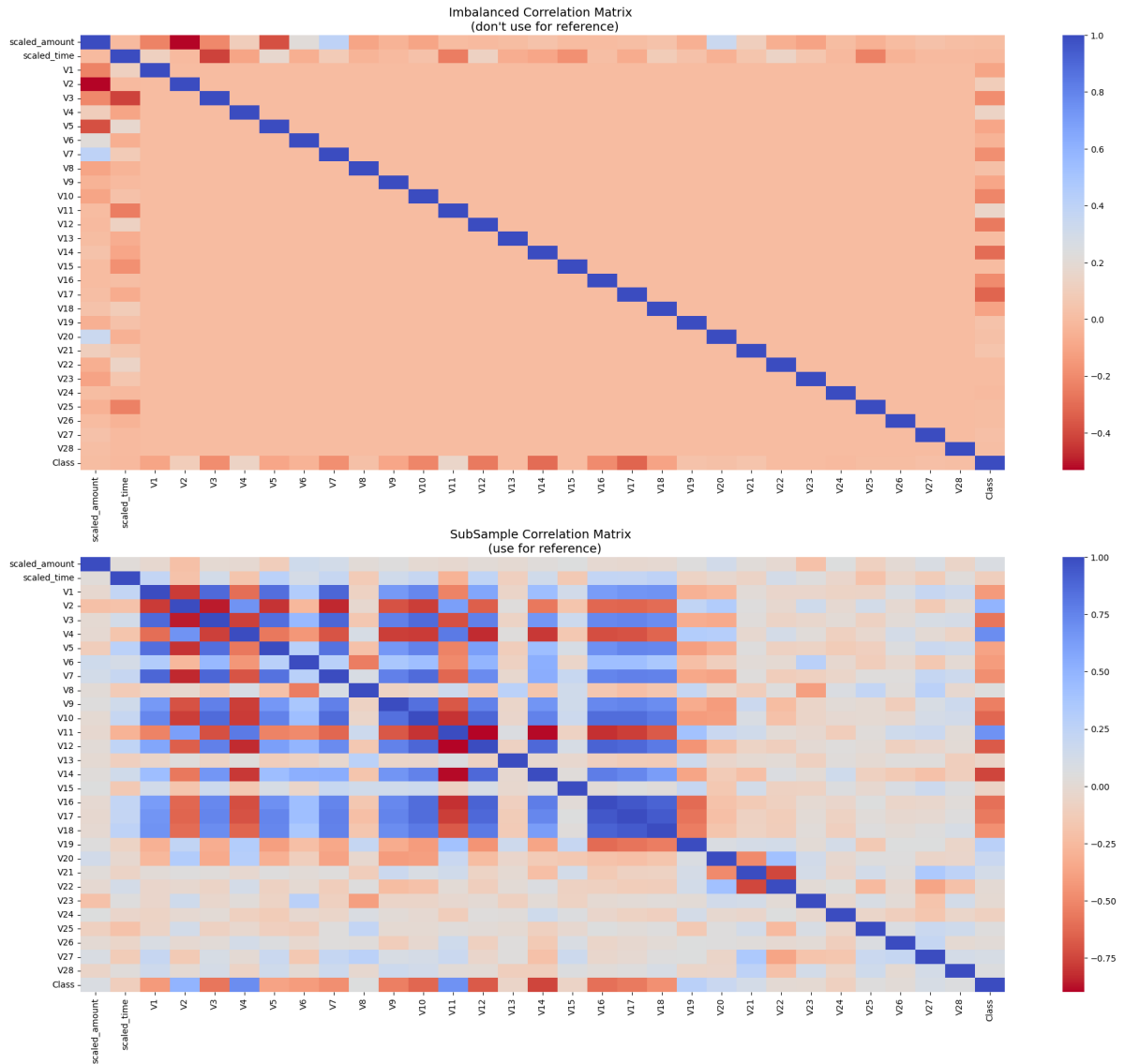
In [ ]:
# Make sure we use the subsample in our correlation

f, (ax1, ax2) = plt.subplots(2, 1, figsize=(24,20))

# Entire DataFrame
corr = df.corr()
sns.heatmap(corr, cmap='coolwarm_r', annot_kws={'size':20}, ax=ax1)
ax1.set_title("Imbalanced Correlation Matrix \n (don't use for reference)", fontsize=14)

```

```
sub_sample_corr = new_df.corr()
sns.heatmap(sub_sample_corr, cmap='coolwarm_r', annot_kws={'size':20}, ax=ax2)
ax2.set_title('SubSample Correlation Matrix \n (use for reference)', fontsize=14)
plt.show()
```



```
In [ ]: f, axes = plt.subplots(ncols=4, figsize=(20,4))

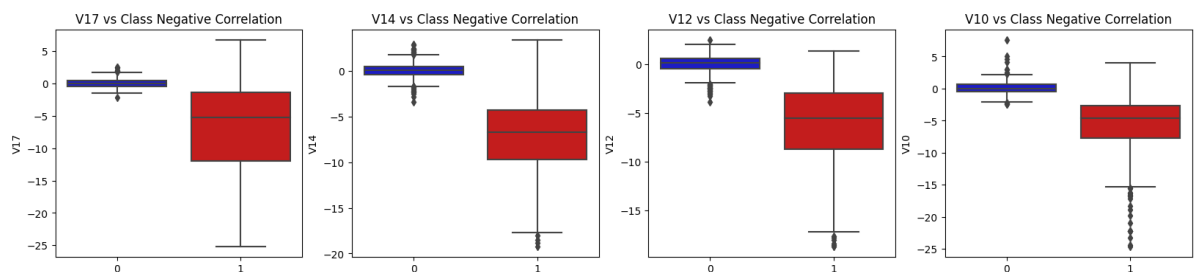
# Negative Correlations with our Class (The lower our feature value the more likely it will be a frau
sns.boxplot(x="Class", y="V17", data=new_df, palette=colors, ax=axes[0])
axes[0].set_title('V17 vs Class Negative Correlation')

sns.boxplot(x="Class", y="V14", data=new_df, palette=colors, ax=axes[1])
axes[1].set_title('V14 vs Class Negative Correlation')

sns.boxplot(x="Class", y="V12", data=new_df, palette=colors, ax=axes[2])
axes[2].set_title('V12 vs Class Negative Correlation')

sns.boxplot(x="Class", y="V10", data=new_df, palette=colors, ax=axes[3])
axes[3].set_title('V10 vs Class Negative Correlation')

plt.show()
```



```
In [ ]: f, axes = plt.subplots(ncols=4, figsize=(20,4))

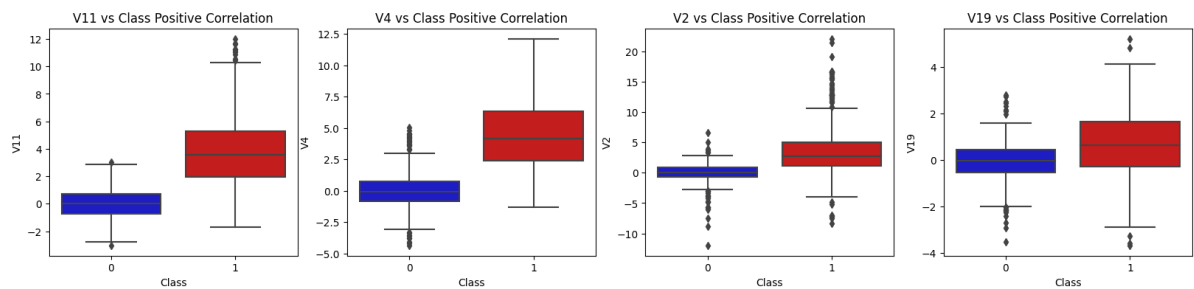
# Positive correlations (The higher the feature the probability increases that it will be a fraud tra
sns.boxplot(x="Class", y="V11", data=new_df, palette=colors, ax=axes[0])
axes[0].set_title('V11 vs Class Positive Correlation')

sns.boxplot(x="Class", y="V4", data=new_df, palette=colors, ax=axes[1])
axes[1].set_title('V4 vs Class Positive Correlation')

sns.boxplot(x="Class", y="V2", data=new_df, palette=colors, ax=axes[2])
axes[2].set_title('V2 vs Class Positive Correlation')

sns.boxplot(x="Class", y="V19", data=new_df, palette=colors, ax=axes[3])
axes[3].set_title('V19 vs Class Positive Correlation')

plt.show()
```



```
In [ ]: # histogram
from scipy.stats import norm

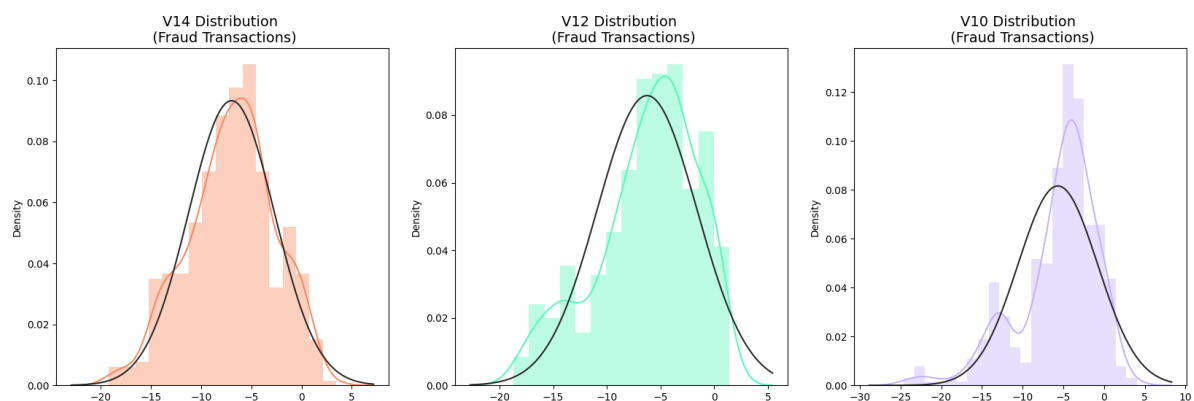
f, (ax1, ax2, ax3) = plt.subplots(1,3, figsize=(20, 6))

v14_fraud_dist = new_df['V14'].loc[new_df['Class'] == 1].values
sns.distplot(v14_fraud_dist,ax=ax1, fit=norm, color='#FB8861')
ax1.set_title('V14 Distribution \n (Fraud Transactions)', fontsize=14)

v12_fraud_dist = new_df['V12'].loc[new_df['Class'] == 1].values
sns.distplot(v12_fraud_dist,ax=ax2, fit=norm, color='#56F9BB')
ax2.set_title('V12 Distribution \n (Fraud Transactions)', fontsize=14)

v10_fraud_dist = new_df['V10'].loc[new_df['Class'] == 1].values
sns.distplot(v10_fraud_dist,ax=ax3, fit=norm, color='#C5B3F9')
ax3.set_title('V10 Distribution \n (Fraud Transactions)', fontsize=14)

plt.show()
```



```
In [ ]: f,(ax1, ax2, ax3) = plt.subplots(1, 3, figsize=(20,6))

colors = ['#B3F9C5', '#f9c5b3']
# Boxplots with outliers removed
# Feature V14
sns.boxplot(x="Class", y="V14", data=new_df,ax=ax1, palette=colors)
ax1.set_title("V14 Feature \n Reduction of outliers", fontsize=14)
ax1.annotate('Fewer extreme \n outliers', xy=(0.98, -17.5), xytext=(0, -12),
            arrowprops=dict(facecolor='black'),
```

```

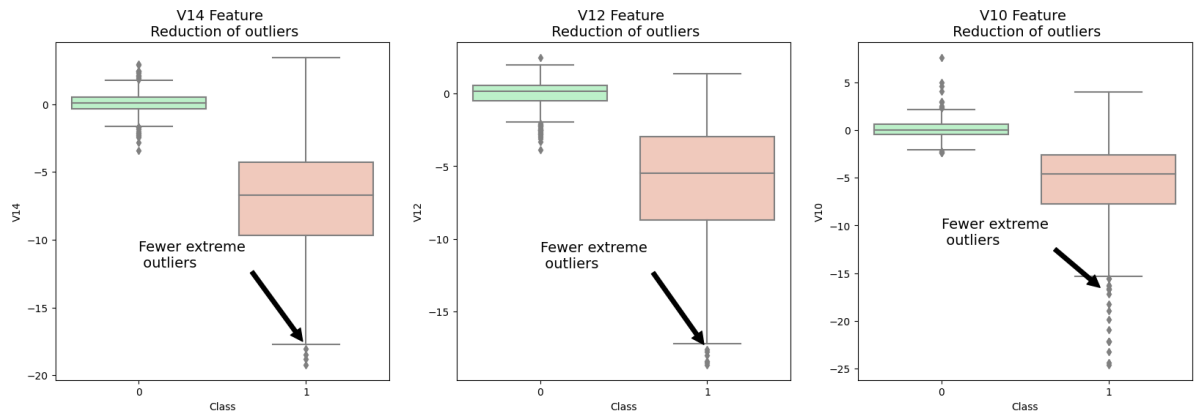
        fontsize=14)

# Feature 12
sns.boxplot(x="Class", y="V12", data=new_df, ax=ax2, palette=colors)
ax2.set_title("V12 Feature \n Reduction of outliers", fontsize=14)
ax2.annotate('Fewer extreme \n outliers', xy=(0.98, -17.3), xytext=(0, -12),
            arrowprops=dict(facecolor='black'),
            fontsize=14)

# Feature V10
sns.boxplot(x="Class", y="V10", data=new_df, ax=ax3, palette=colors)
ax3.set_title("V10 Feature \n Reduction of outliers", fontsize=14)
ax3.annotate('Fewer extreme \n outliers', xy=(0.95, -16.5), xytext=(0, -12),
            arrowprops=dict(facecolor='black'),
            fontsize=14)

plt.show()

```



```

In [ ]: # implementation timings
X = new_df.drop('Class', axis=1)
y = new_df['Class']

# T-SNE Implementation
t0 = time.time()
X_reduced_tsne = TSNE(n_components=2, random_state=42).fit_transform(X.values)
t1 = time.time()
print("T-SNE took {:.2} s".format(t1 - t0))

# PCA Implementation
t0 = time.time()
X_reduced_pca = PCA(n_components=2, random_state=42).fit_transform(X.values)
t1 = time.time()
print("PCA took {:.2} s".format(t1 - t0))

# TruncatedSVD
t0 = time.time()
X_reduced_svd = TruncatedSVD(n_components=2, algorithm='randomized', random_state=42).fit_transform(X.values)
t1 = time.time()
print("Truncated SVD took {:.2} s".format(t1 - t0))

T-SNE took 3.7 s
PCA took 0.0065 s
Truncated SVD took 0.005 s

```

```

In [ ]: #clustering
f, (ax1, ax2, ax3) = plt.subplots(1, 3, figsize=(24,6))

```

```

import matplotlib.pyplot as plt
import seaborn as sns
from sklearn.manifold import TSNE
from sklearn.decomposition import PCA, TruncatedSVD
import matplotlib.patches as mpatches
import time

# Classifier Libraries
from sklearn.linear_model import LogisticRegression
from sklearn.svm import SVC
from sklearn.neighbors import KNeighborsClassifier
from sklearn.tree import DecisionTreeClassifier
from sklearn.ensemble import RandomForestClassifier
import collections

```

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

In [3]:

```

Total_transactions = len(df)
normal = len(df[df.Class == 0])
fraudulent = len(df[df.Class == 1])
fraud_percentage = round(fraudulent/normal*100, 2)

print("total number of transactions: ",Total_transactions)
print("normal transactions: ",normal)
print("fraud transactions: ",fraudulent)
print("percentage: ",fraud_percentage)

```

```

total number of transactions: 284807
normal transactions: 284315
fraud transactions: 492
percentage: 0.17

```

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

Out[4]:

	Time	V1	V2	V3	V4	V5	V6
0	0.0	-1.359807	-0.072781	2.536347	1.378155	-0.338321	0.4623
1	0.0	1.191857	0.266151	0.166480	0.448154	0.060018	-0.0823
2	1.0	-1.358354	-1.340163	1.773209	0.379780	-0.503198	1.8004
3	1.0	-0.966272	-0.185226	1.792993	-0.863291	-0.010309	1.2472
4	2.0	-1.158233	0.877737	1.548718	0.403034	-0.407193	0.0959

5 rows × 31 columns

In [5]: `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
```

```
In [6]: print(df.duplicated().sum())
df.shape
```

1081

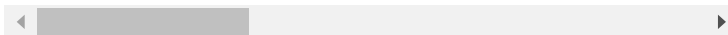
Out[6]: (284807, 31)

```
In [7]: df.drop_duplicates()
```

Out[7]:

	Time	V1	V2	V3	V4	
0	0.0	-1.359807	-0.072781	2.536347	1.378155	-0.3383
1	0.0	1.191857	0.266151	0.166480	0.448154	0.0600
2	1.0	-1.358354	-1.340163	1.773209	0.379780	-0.5031
3	1.0	-0.966272	-0.185226	1.792993	-0.863291	-0.0103
4	2.0	-1.158233	0.877737	1.548718	0.403034	-0.4071
...
284802	172786.0	-11.881118	10.071785	-9.834783	-2.066656	-5.3644
284803	172787.0	-0.732789	-0.055080	2.035030	-0.738589	0.8682
284804	172788.0	1.919565	-0.301254	-3.249640	-0.557828	2.6305
284805	172788.0	-0.240440	0.530483	0.702510	0.689799	-0.3779
284806	172792.0	-0.533413	-0.189733	0.703337	-0.506271	-0.0125

283726 rows × 31 columns

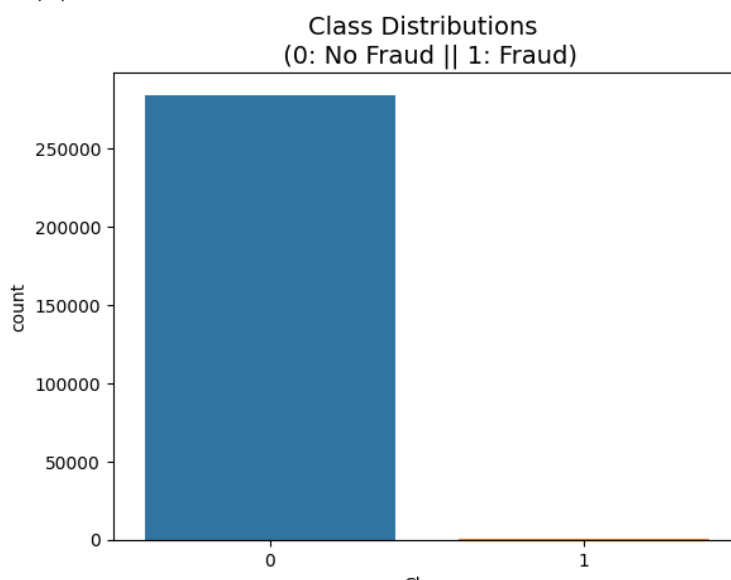


Pre-Processing

```
In [8]: import matplotlib.pyplot as plt
colors = ["#0101DF", "#DF0101"]

sns.countplot(x='Class', data=df)
plt.title('Class Distributions \n (0: No Fraud || 1: Fraud)', f
```

Out[8]: Text(0.5, 1.0, 'Class Distributions \n (0: No Fraud || 1: Fraud)')




```
In [9]: import warnings
warnings.filterwarnings("ignore")
```

```
In [10]: fig, ax = plt.subplots(1, 2, figsize=(18,4))

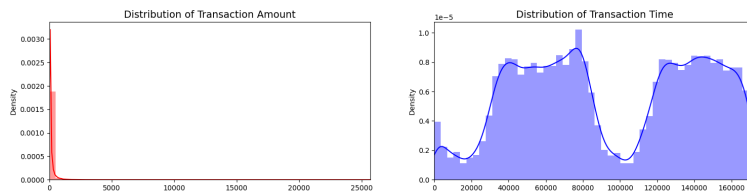
amount_val = df['Amount'].values
time_val = df['Time'].values

sns.distplot(amount_val, ax=ax[0], color='r')
ax[0].set_title('Distribution of Transaction Amount', fontsize=14)
ax[0].set_xlim([min(amount_val), max(amount_val)])

sns.distplot(time_val, ax=ax[1], color='b')
ax[1].set_title('Distribution of Transaction Time', fontsize=14)
ax[1].set_xlim([min(time_val), max(time_val)])

plt.show()

#very imbalanced data in amount and time
```



```
In [11]: #scaling columns amount and time

from sklearn.preprocessing import StandardScaler, RobustScaler

# RobustScaler

std_scaler = StandardScaler()
rob_scaler = RobustScaler()

df['scaled_amount'] = rob_scaler.fit_transform(df['Amount'].values)
df['scaled_time'] = rob_scaler.fit_transform(df['Time'].values)

df.drop(['Time', 'Amount'], axis=1, inplace=True)
```

```
In [12]: scaled_amount = df['scaled_amount']
scaled_time = df['scaled_time']

df.drop(['scaled_amount', 'scaled_time'], axis=1, inplace=True)
df.insert(0, 'scaled_amount', scaled_amount)
df.insert(1, 'scaled_time', scaled_time)

# Amount and Time are Scaled!

df.head()
```

```
Out[12]:
```

	scaled_amount	scaled_time	V1	V2	V3	
0	1.783274	-0.994983	-1.359807	-0.072781	2.536347	1.378
1	-0.269825	-0.994983	1.191857	0.266151	0.166480	0.448
2	4.983721	-0.994972	-1.358354	-1.340163	1.773209	0.379
3	1.418291	-0.994972	-0.966272	-0.185226	1.792993	-0.863
4	0.670579	-0.994960	-1.158233	0.877737	1.548718	0.403

5 rows × 31 columns

```
In [13]: fig, ax = plt.subplots(1, 2, figsize=(18,4))
```

```

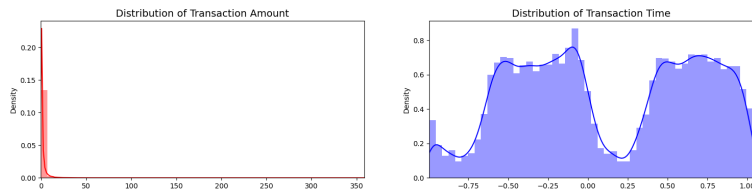
amount_val = df['scaled_amount'].values
time_val = df['scaled_time'].values

sns.distplot(amount_val, ax=ax[0], color='r')
ax[0].set_title('Distribution of Transaction Amount', fontsize=14)
ax[0].set_xlim([min(amount_val), max(amount_val)])

sns.distplot(time_val, ax=ax[1], color='b')
ax[1].set_title('Distribution of Transaction Time', fontsize=14)
ax[1].set_xlim([min(time_val), max(time_val)])

plt.show()

```



```

In [14]: from sklearn.model_selection import train_test_split
from sklearn.model_selection import StratifiedShuffleSplit
from sklearn.model_selection import KFold, StratifiedKFold

print('No Frauds', round(df['Class'].value_counts()[0]/len(df)
print('Frauds', round(df['Class'].value_counts()[1]/len(df) * 100)

X = df.drop('Class', axis=1)
y = df['Class']

sss = StratifiedKFold(n_splits=5, random_state=None, shuffle=False)

for train_index, test_index in sss.split(X, y):
    print("Train:", train_index, "Test:", test_index)
    original_Xtrain, original_Xtest = X.iloc[train_index], X.iloc[test_index]
    original_ytrain, original_ytest = y.iloc[train_index], y.iloc[test_index]

```

```

No Frauds 99.83 % of the dataset
Frauds 0.17 % of the dataset
Train: [ 30473  30496  31002 ... 284804 284805 284806] Test: [
0      1      2 ... 57017 57018 57019]
Train: [      0      1      2 ... 284804 284805 284806] Test: [ 30
473  30496  31002 ... 113964 113965 113966]
Train: [      0      1      2 ... 284804 284805 284806] Test: [ 81
609  82400  83053 ... 170946 170947 170948]
Train: [      0      1      2 ... 284804 284805 284806] Test: [150
654 150660 150661 ... 227866 227867 227868]
Train: [      0      1      2 ... 227866 227867 227868] Test: [212
516 212644 213092 ... 284804 284805 284806]

```

```

In [15]: # Random Under Sampling" which basically consists of removing a
# Lets shuffle the data before creating the subsamples

df = df.sample(frac=1)

# amount of fraud classes 492 rows.
fraud_df = df.loc[df['Class'] == 1]
non_fraud_df = df.loc[df['Class'] == 0][:492]

normal_distributed_df = pd.concat([fraud_df, non_fraud_df])

# Shuffle dataframe rows
new_df = normal_distributed_df.sample(frac=1, random_state=42)

new_df.head()

```

```

Out[15]:

```

	scaled_amount	scaled_time	V1	V2	V3
153730	-0.027947	0.179114	-2.004446	0.990127	2.554818
63634	1.089779	-0.398078	-9.169790	7.092197	-12.354037
92943	0.070286	-0.240980	1.079165	-0.590200	1.936351
12108	-0.293440	-0.747730	-16.917468	9.669900	-23.736447

244004 -0.293440 0.792690 -4.673231 4.195976 -8.392423

5 rows × 31 columns

```
In [16]: print('Distribution of the Classes in the subsample dataset')
print(new_df['Class'].value_counts()/len(new_df))

sns.countplot(x='Class', data=new_df, palette=colors)
plt.title('Equally Distributed Classes', fontsize=14)
plt.show()
```

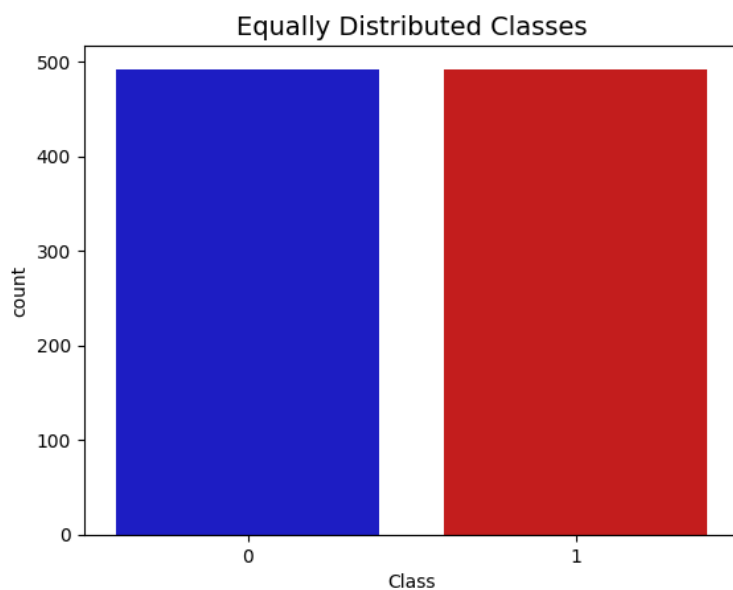
Distribution of the Classes in the subsample dataset

Class

0 0.5

1 0.5

Name: count, dtype: float64

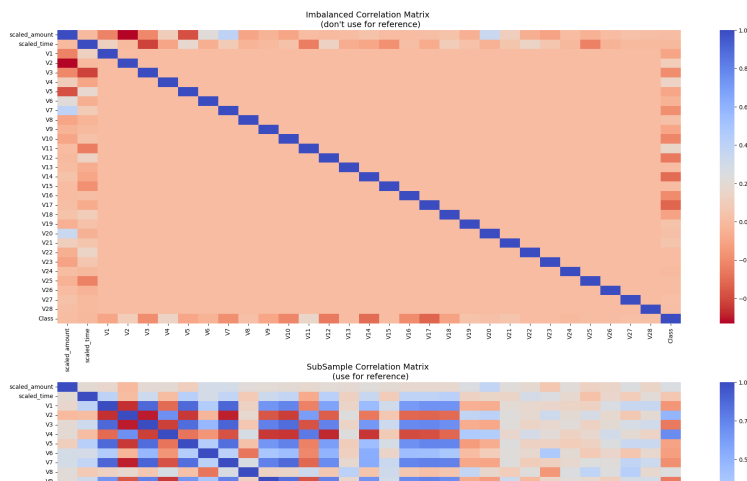


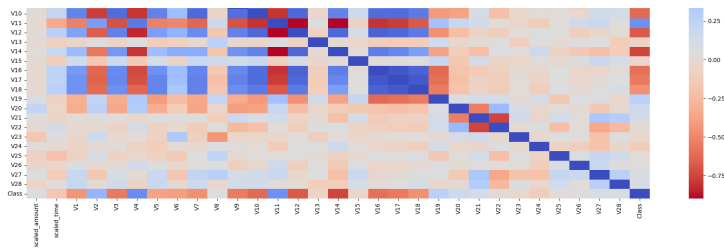
```
In [17]: # Make sure we use the subsample in our correlation

f, (ax1, ax2) = plt.subplots(2, 1, figsize=(24,20))

# Entire DataFrame
corr = df.corr()
sns.heatmap(corr, cmap='coolwarm_r', annot_kws={'size':20}, ax=
ax1.set_title("Imbalanced Correlation Matrix \n (don't use for

sub_sample_corr = new_df.corr()
sns.heatmap(sub_sample_corr, cmap='coolwarm_r', annot_kws={'siz
ax2.set_title('SubSample Correlation Matrix \n (use for referer
plt.show()
```





```
In [18]: f, axes = plt.subplots(ncols=4, figsize=(20,4))

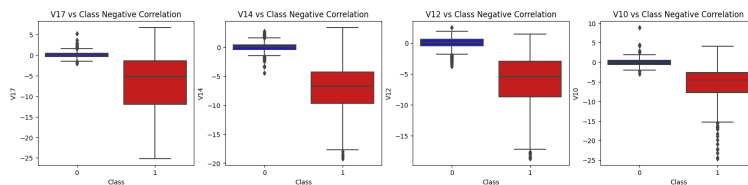
# Negative Correlations with our Class (The lower our feature value the higher the probability of being in Class 1)
sns.boxplot(x="Class", y="V17", data=new_df, palette=colors, ax=axes[0].set_title('V17 vs Class Negative Correlation'))

sns.boxplot(x="Class", y="V14", data=new_df, palette=colors, ax=axes[1].set_title('V14 vs Class Negative Correlation'))

sns.boxplot(x="Class", y="V12", data=new_df, palette=colors, ax=axes[2].set_title('V12 vs Class Negative Correlation'))

sns.boxplot(x="Class", y="V10", data=new_df, palette=colors, ax=axes[3].set_title('V10 vs Class Negative Correlation'))

plt.show()
```



```
In [19]: f, axes = plt.subplots(ncols=4, figsize=(20,4))

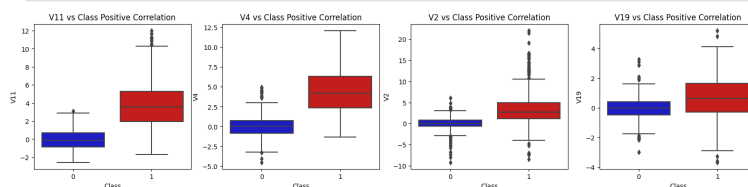
# Positive correlations (The higher the feature the probability of being in Class 1)
sns.boxplot(x="Class", y="V11", data=new_df, palette=colors, ax=axes[0].set_title('V11 vs Class Positive Correlation'))

sns.boxplot(x="Class", y="V4", data=new_df, palette=colors, ax=axes[1].set_title('V4 vs Class Positive Correlation'))

sns.boxplot(x="Class", y="V2", data=new_df, palette=colors, ax=axes[2].set_title('V2 vs Class Positive Correlation'))

sns.boxplot(x="Class", y="V19", data=new_df, palette=colors, ax=axes[3].set_title('V19 vs Class Positive Correlation'))

plt.show()
```



```
In [20]: from scipy.stats import norm

f, (ax1, ax2, ax3) = plt.subplots(1,3, figsize=(20, 6))

v14_fraud_dist = new_df['V14'].loc[new_df['Class'] == 1].values
sns.distplot(v14_fraud_dist,ax=ax1, fit=norm, color='#FB8861')
ax1.set_title('V14 Distribution \n (Fraud Transactions)', fontstyle='italic')

v12_fraud_dist = new_df['V12'].loc[new_df['Class'] == 1].values
sns.distplot(v12_fraud_dist,ax=ax2, fit=norm, color='#56F9BB')
ax2.set_title('V12 Distribution \n (Fraud Transactions)', fontstyle='italic')

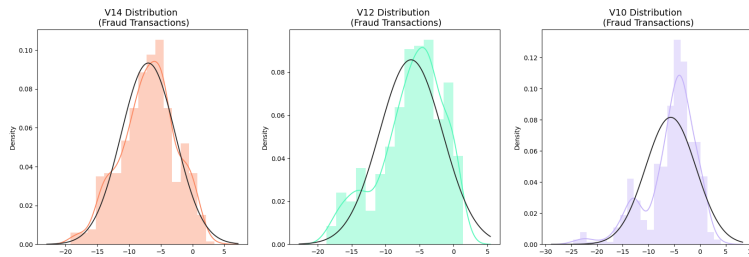
v10_fraud_dist = new_df['V10'].loc[new_df['Class'] == 1].values
```

```

sns.distplot(v10_fraud_dist, ax=ax3, fit=norm, color='#C5B3F9')
ax3.set_title('V10 Distribution \n (Fraud Transactions)', fontsize=14)

plt.show()

```



```

In [21]: f,(ax1, ax2, ax3) = plt.subplots(1, 3, figsize=(20,6))

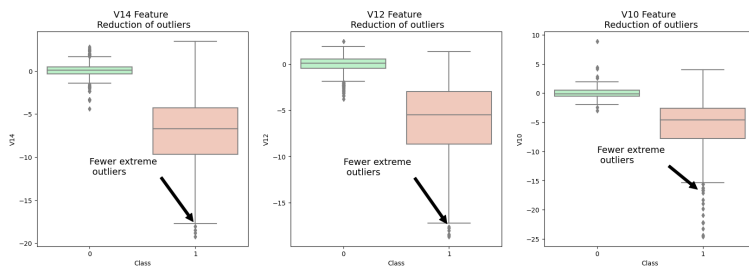
colors = ['#B3F9C5', '#f9c5b3']
# Boxplots with outliers removed
# Feature V14
sns.boxplot(x="Class", y="V14", data=new_df, ax=ax1, palette=colors)
ax1.set_title("V14 Feature \n Reduction of outliers", fontsize=14)
ax1.annotate('Fewer extreme \n outliers', xy=(0.98, -17.5), xytext=(0.9, -18.5),
            arrowprops=dict(facecolor='black'),
            fontsize=14)

# Feature 12
sns.boxplot(x="Class", y="V12", data=new_df, ax=ax2, palette=colors)
ax2.set_title("V12 Feature \n Reduction of outliers", fontsize=14)
ax2.annotate('Fewer extreme \n outliers', xy=(0.98, -17.3), xytext=(0.9, -18.3),
            arrowprops=dict(facecolor='black'),
            fontsize=14)

# Feature V10
sns.boxplot(x="Class", y="V10", data=new_df, ax=ax3, palette=colors)
ax3.set_title("V10 Feature \n Reduction of outliers", fontsize=14)
ax3.annotate('Fewer extreme \n outliers', xy=(0.95, -16.5), xytext=(0.9, -17.5),
            arrowprops=dict(facecolor='black'),
            fontsize=14)

plt.show()

```



```

In [22]: X = new_df.drop('Class', axis=1)
y = new_df['Class']

# T-SNE Implementation
t0 = time.time()
X_reduced_tsne = TSNE(n_components=2, random_state=42).fit_transform(X)
t1 = time.time()
print("T-SNE took {:.2} s".format(t1 - t0))

# PCA Implementation
t0 = time.time()
X_reduced_pca = PCA(n_components=2, random_state=42).fit_transform(X)
t1 = time.time()
print("PCA took {:.2} s".format(t1 - t0))

# TruncatedSVD
t0 = time.time()
X_reduced_svd = TruncatedSVD(n_components=2, algorithm='randomized').fit_transform(X)
t1 = time.time()
print("Truncated SVD took {:.2} s".format(t1 - t0))

```

T-SNE took 4.6 s

PCA took 0.0071 s
Truncated SVD took 0.0066 s

```
In [23]: f, (ax1, ax2, ax3) = plt.subplots(1, 3, figsize=(24,6))
# Labels = ['No Fraud', 'Fraud']
f.suptitle('Clusters using Dimensionality Reduction', fontsize=14)

blue_patch = mpatches.Patch(color='#0A0AFF', label='No Fraud')
red_patch = mpatches.Patch(color='#AF0000', label='Fraud')

# t-SNE scatter plot
ax1.scatter(X_reduced_tsne[:,0], X_reduced_tsne[:,1], c=(y == 0))
ax1.scatter(X_reduced_tsne[:,0], X_reduced_tsne[:,1], c=(y == 1))
ax1.set_title('t-SNE', fontsize=14)

ax1.grid(True)

ax1.legend(handles=[blue_patch, red_patch])

# PCA scatter plot
ax2.scatter(X_reduced_pca[:,0], X_reduced_pca[:,1], c=(y == 0))
ax2.scatter(X_reduced_pca[:,0], X_reduced_pca[:,1], c=(y == 1))
ax2.set_title('PCA', fontsize=14)

ax2.grid(True)

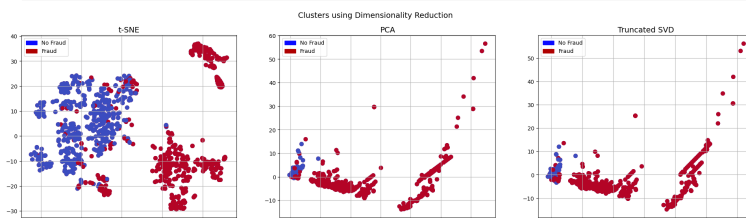
ax2.legend(handles=[blue_patch, red_patch])

# TruncatedSVD scatter plot
ax3.scatter(X_reduced_svd[:,0], X_reduced_svd[:,1], c=(y == 0))
ax3.scatter(X_reduced_svd[:,0], X_reduced_svd[:,1], c=(y == 1))
ax3.set_title('Truncated SVD', fontsize=14)

ax3.grid(True)

ax3.legend(handles=[blue_patch, red_patch])

plt.show()
```



Classification Models

```
In [24]: # undersampling
X = new_df.drop('Class', axis=1)
y = new_df['Class']
```

```
In [25]: from sklearn.model_selection import train_test_split
# splitting
X_train, X_test, y_train, y_test = train_test_split(X, y, test_
X_train = X_train.values
X_test = X_test.values
y_train = y_train.values
y_test = y_test.values
```

```
In [18]: import numpy as np # linear algebra
import pandas as pd # data processing, CSV file I/O (e.g. pd.read_csv)
import tensorflow as tf
import matplotlib.pyplot as plt
import seaborn as sns
from sklearn.manifold import TSNE
from sklearn.decomposition import PCA, TruncatedSVD
import matplotlib.patches as mpatches
import time

# Classifier Libraries
from sklearn.linear_model import LogisticRegression
from sklearn.svm import SVC
from sklearn.neighbors import KNeighborsClassifier
from sklearn.tree import DecisionTreeClassifier
from sklearn.ensemble import RandomForestClassifier
import collections
```

```
In [19]: df = pd.read_csv("creditcard.csv")
```

```
In [20]: Total_transactions = len(df)
normal = len(df[df.Class == 0])
fraudulent = len(df[df.Class == 1])
fraud_percentage = round(fraudulent/normal*100, 2)

print("total number of transactions: ",Total_transactions)
print("normal transactions: ",normal)
print("fraud transactions: ",fraudulent)
print("percentage: ",fraud_percentage)
```

```
total number of transactions: 284807
normal transactions: 284315
fraud transactions: 492
percentage: 0.17
```

```
In [21]: df.head()
```

```
Out[21]:
```

	Time	V1	V2	V3	V4	V5	V6	V7	V8	V9	...
0	0.0	-1.359807	-0.072781	2.536347	1.378155	-0.338321	0.462388	0.239599	0.098698	0.363787	...
1	0.0	1.191857	0.266151	0.166480	0.448154	0.060018	-0.082361	-0.078803	0.085102	-0.255425	...
2	1.0	-1.358354	-1.340163	1.773209	0.379780	-0.503198	1.800499	0.791461	0.247676	-1.514654	...
3	1.0	-0.966272	-0.185226	1.792993	-0.863291	-0.010309	1.247203	0.237609	0.377436	-1.387024	...
4	2.0	-1.158233	0.877737	1.548718	0.403034	-0.407193	0.095921	0.592941	-0.270533	0.817739	...

5 rows × 31 columns

```
In [22]: df.isnull().sum()
```

```
Out[22]: 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 [23]: print(df.duplicated().sum())
df.shape
```

1081

```
Out[23]: (284807, 31)
```

```
In [24]: df.drop_duplicates()
```

```
Out[24]:
```

	Time	V1	V2	V3	V4	V5	V6	V7	V8	
0	0.0	-1.359807	-0.072781	2.536347	1.378155	-0.338321	0.462388	0.239599	0.098698	0.
1	0.0	1.191857	0.266151	0.166480	0.448154	0.060018	-0.082361	-0.078803	0.085102	-0.
2	1.0	-1.358354	-1.340163	1.773209	0.379780	-0.503198	1.800499	0.791461	0.247676	-1.
3	1.0	-0.966272	-0.185226	1.792993	-0.863291	-0.010309	1.247203	0.237609	0.377436	-1.
4	2.0	-1.158233	0.877737	1.548718	0.403034	-0.407193	0.095921	0.592941	-0.270533	0.
...
284802	172786.0	-11.881118	10.071785	-9.834783	-2.066656	-5.364473	-2.606837	-4.918215	7.305334	1.
284803	172787.0	-0.732789	-0.055080	2.035030	-0.738589	0.868229	1.058415	0.024330	0.294869	0.
284804	172788.0	1.919565	-0.301254	-3.249640	-0.557828	2.630515	3.031260	-0.296827	0.708417	0.
284805	172788.0	-0.240440	0.530483	0.702510	0.689799	-0.377961	0.623708	-0.686180	0.679145	0.
284806	172792.0	-0.533413	-0.189733	0.703337	-0.506271	-0.012546	-0.649617	1.577006	-0.414650	0.

283726 rows × 31 columns

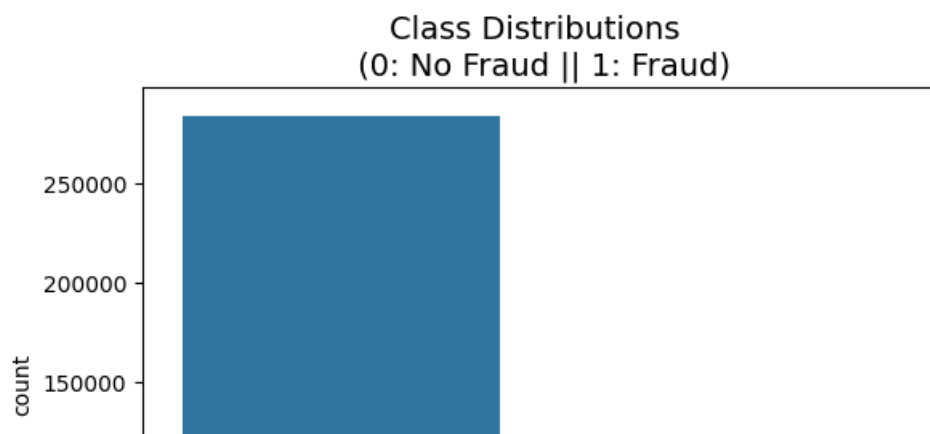


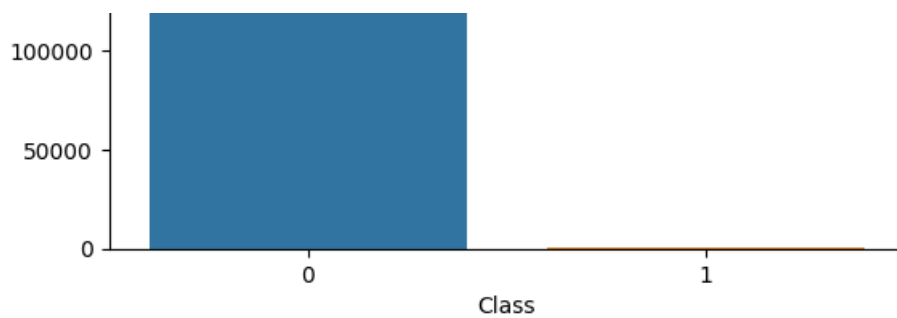
Pre-Processing

```
In [25]: import matplotlib.pyplot as plt
colors = ["#0101DF", "#DF0101"]

sns.countplot(x='Class', data=df)
plt.title('Class Distributions \n (0: No Fraud || 1: Fraud)', fontsize=14)
```

```
Out[25]: Text(0.5, 1.0, 'Class Distributions \n (0: No Fraud || 1: Fraud)')
```





```
In [26]: import warnings
warnings.filterwarnings("ignore")
```

```
In [27]: fig, ax = plt.subplots(1, 2, figsize=(18,4))

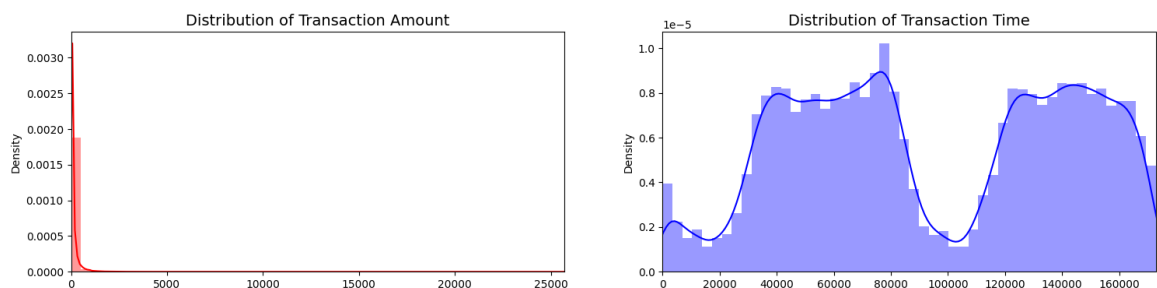
amount_val = df['Amount'].values
time_val = df['Time'].values

sns.distplot(amount_val, ax=ax[0], color='r')
ax[0].set_title('Distribution of Transaction Amount', fontsize=14)
ax[0].set_xlim([min(amount_val), max(amount_val)])

sns.distplot(time_val, ax=ax[1], color='b')
ax[1].set_title('Distribution of Transaction Time', fontsize=14)
ax[1].set_xlim([min(time_val), max(time_val)])

plt.show()

#very imbalanced data in amount and time
```



```
In [28]: #scaling columns amount and time

from sklearn.preprocessing import StandardScaler, RobustScaler

# RobustScaler

std_scaler = StandardScaler()
rob_scaler = RobustScaler()

df['scaled_amount'] = rob_scaler.fit_transform(df['Amount'].values.reshape(-1,1))
df['scaled_time'] = rob_scaler.fit_transform(df['Time'].values.reshape(-1,1))

df.drop(['Time', 'Amount'], axis=1, inplace=True)
```

```
In [29]: scaled_amount = df['scaled_amount']
scaled_time = df['scaled_time']

df.drop(['scaled_amount', 'scaled_time'], axis=1, inplace=True)
df.insert(0, 'scaled_amount', scaled_amount)
df.insert(1, 'scaled_time', scaled_time)

# Amount and Time are Scaled!

df.head()
```

```
Out[29]:
```

	scaled_amount	scaled_time	V1	V2	V3	V4	V5	V6	V7
0	1.783274	-0.994983	-1.359807	-0.072781	2.536347	1.378155	-0.338321	0.462388	0.239599

1	-0.269825	-0.994983	1.191857	0.266151	0.166480	0.448154	0.060018	-0.082361	-0.078803	0
2	4.983721	-0.994972	-1.358354	-1.340163	1.773209	0.379780	-0.503198	1.800499	0.791461	0
3	1.418291	-0.994972	-0.966272	-0.185226	1.792993	-0.863291	-0.010309	1.247203	0.237609	0
4	0.670579	-0.994960	-1.158233	0.877737	1.548718	0.403034	-0.407193	0.095921	0.592941	-0

5 rows × 31 columns

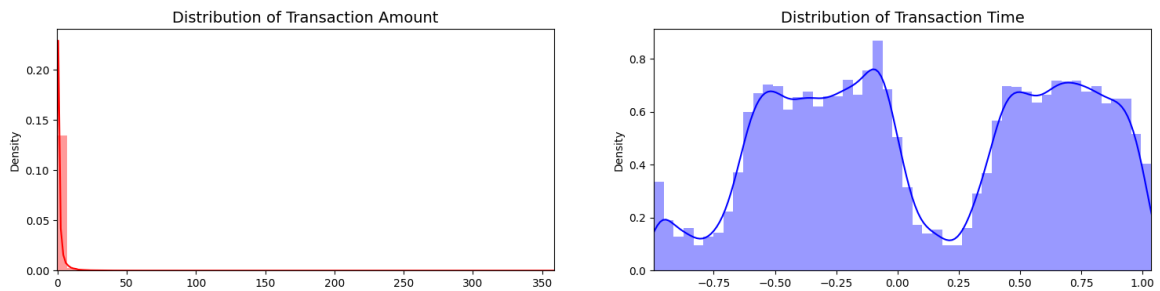
```
In [30]: fig, ax = plt.subplots(1, 2, figsize=(18,4))

amount_val = df['scaled_amount'].values
time_val = df['scaled_time'].values

sns.distplot(amount_val, ax=ax[0], color='r')
ax[0].set_title('Distribution of Transaction Amount', fontsize=14)
ax[0].set_xlim([min(amount_val), max(amount_val)])

sns.distplot(time_val, ax=ax[1], color='b')
ax[1].set_title('Distribution of Transaction Time', fontsize=14)
ax[1].set_xlim([min(time_val), max(time_val)])

plt.show()
```



```
In [31]: from sklearn.model_selection import train_test_split
from sklearn.model_selection import StratifiedShuffleSplit
from sklearn.model_selection import KFold, StratifiedKFold

print('No Frauds', round(df['Class'].value_counts()[0]/len(df) * 100,2), '% of the dataset')
print('Frauds', round(df['Class'].value_counts()[1]/len(df) * 100,2), '% of the dataset')

X = df.drop('Class', axis=1)
y = df['Class']

sss = StratifiedKFold(n_splits=5, random_state=None, shuffle=False)

for train_index, test_index in sss.split(X, y):
    print("Train:", train_index, "Test:", test_index)
    original_Xtrain, original_Xtest = X.iloc[train_index], X.iloc[test_index]
    original_ytrain, original_ytest = y.iloc[train_index], y.iloc[test_index]
```

```
No Frauds 99.83 % of the dataset
Frauds 0.17 % of the dataset
Train: [ 30473  30496  31002 ... 284804 284805 284806] Test: [    0     1     2 ... 57017 57018 57019]
Train: [    0     1     2 ... 284804 284805 284806] Test: [ 30473  30496  31002 ... 113964 113965 113966]
Train: [    0     1     2 ... 284804 284805 284806] Test: [ 81609  82400  83053 ... 170946 170947 170948]
Train: [    0     1     2 ... 284804 284805 284806] Test: [150654 150660 150661 ... 227866 227867 227868]
Train: [    0     1     2 ... 227866 227867 227868] Test: [212516 212644 213092 ... 284804 284805 284806]
```

```
In [32]: # Random Under Sampling" which basically consists of removing data in order to have a more balance
# Lets shuffle the data before creating the subsamples

df = df.sample(frac=1)

# amount of fraud classes 492 rows.
```

```

fraud_df = df.loc[df['Class'] == 1]
non_fraud_df = df.loc[df['Class'] == 0][:492]

normal_distributed_df = pd.concat([fraud_df, non_fraud_df])

# Shuffle dataframe rows
new_df = normal_distributed_df.sample(frac=1, random_state=42)

new_df.head()

```

Out[32]:

	scaled_amount	scaled_time	V1	V2	V3	V4	V5	V6	
4774	1.229651	-0.945194	-0.519700	1.087853	1.475829	0.595549	0.135176	0.161461	1.
15166	1.089779	-0.683384	-18.474868	11.586381	-21.402917	6.038515	-14.451158	-4.146524	-14.
42553	-0.237686	-0.511554	1.060336	-0.171373	1.350160	1.229241	-0.906235	0.457292	-0.
157918	8.567037	0.304574	-1.101035	-1.674928	-0.573388	5.617556	0.765556	0.440607	1.
280149	0.780968	0.994596	-0.676143	1.126366	-2.213700	0.468308	-1.120541	-0.003346	-2.

5 rows × 31 columns

In [33]:

```

print('Distribution of the Classes in the subsample dataset')
print(new_df['Class'].value_counts()/len(new_df))

sns.countplot(x='Class', data=new_df, palette=colors)
plt.title('Equally Distributed Classes', fontsize=14)
plt.show()

```

Distribution of the Classes in the subsample dataset

Class

0 0.5

1 0.5

Name: count, dtype: float64



In [34]:

```

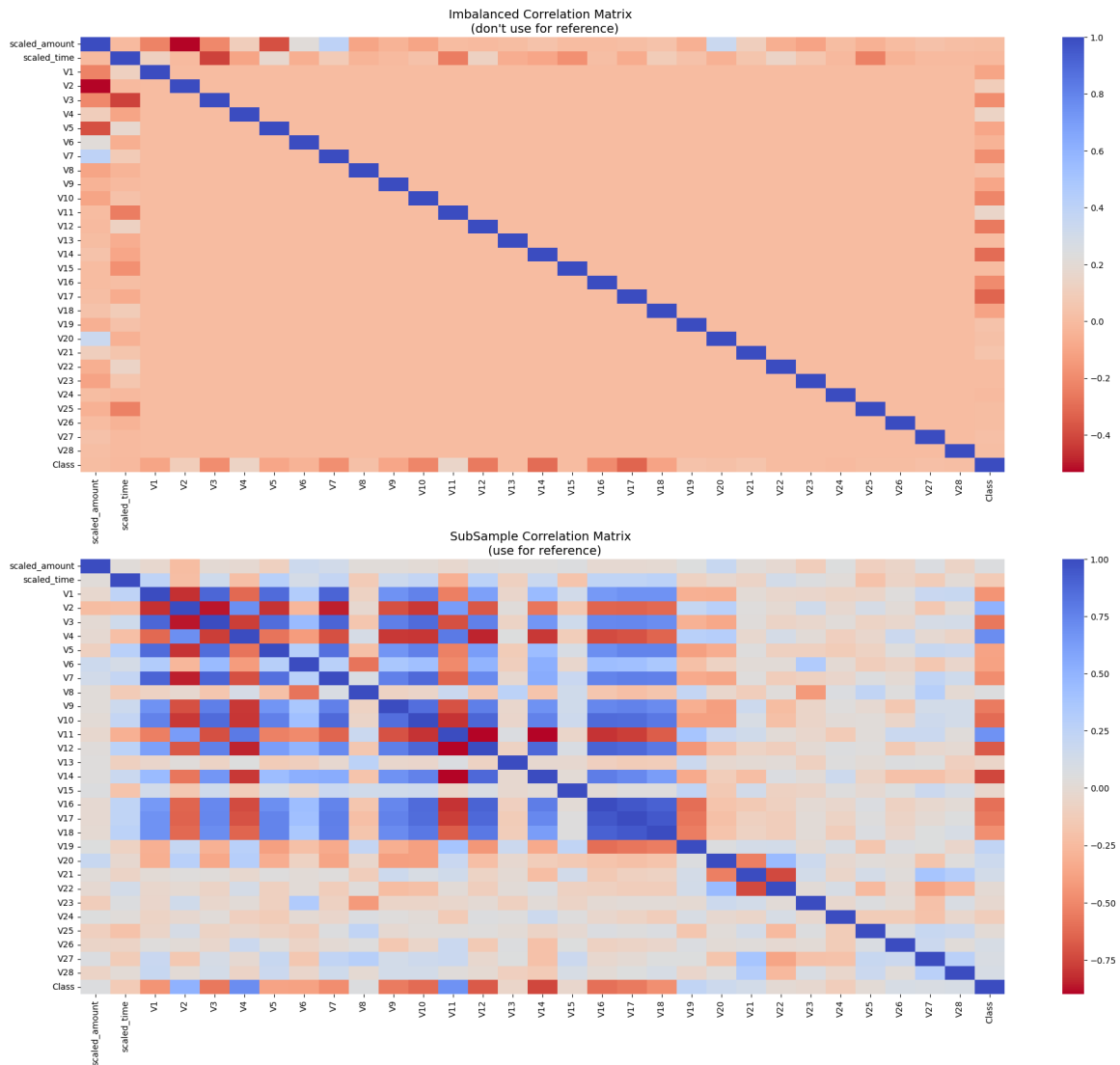
# Make sure we use the subsample in our correlation

f, (ax1, ax2) = plt.subplots(2, 1, figsize=(24,20))

# Entire DataFrame
corr = df.corr()
sns.heatmap(corr, cmap='coolwarm_r', annot_kws={'size':20}, ax=ax1)
ax1.set_title("Imbalanced Correlation Matrix \n (don't use for reference)", fontsize=14)

```

```
sub_sample_corr = new_df.corr()
sns.heatmap(sub_sample_corr, cmap='coolwarm_r', annot_kws={'size':20}, ax=ax2)
ax2.set_title('SubSample Correlation Matrix \n (use for reference)', fontsize=14)
plt.show()
```



```
In [35]: f, axes = plt.subplots(ncols=4, figsize=(20,4))

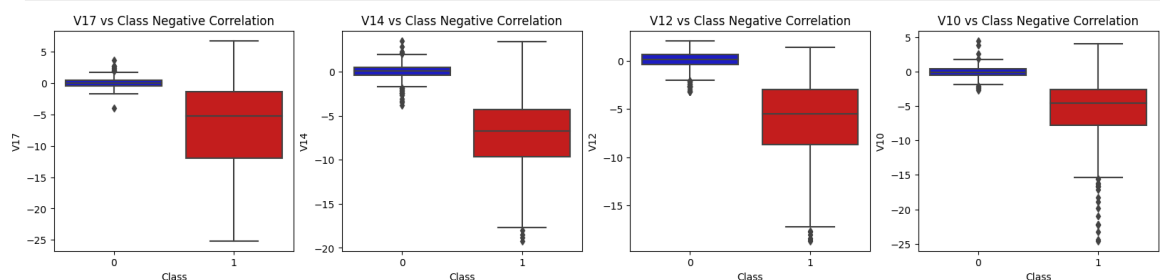
# Negative Correlations with our Class (The lower our feature value the more likely it will be a f
sns.boxplot(x="Class", y="V17", data=new_df, palette=colors, ax=axes[0])
axes[0].set_title('V17 vs Class Negative Correlation')

sns.boxplot(x="Class", y="V14", data=new_df, palette=colors, ax=axes[1])
axes[1].set_title('V14 vs Class Negative Correlation')

sns.boxplot(x="Class", y="V12", data=new_df, palette=colors, ax=axes[2])
axes[2].set_title('V12 vs Class Negative Correlation')

sns.boxplot(x="Class", y="V10", data=new_df, palette=colors, ax=axes[3])
axes[3].set_title('V10 vs Class Negative Correlation')

plt.show()
```



```
In [36]: f, axes = plt.subplots(ncols=4, figsize=(20,4))

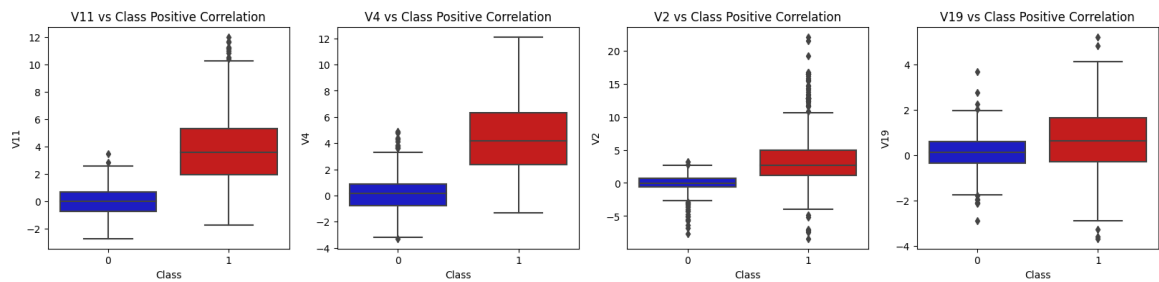
# Positive correlations (The higher the feature the probability increases that it will be a fraud
sns.boxplot(x="Class", y="V11", data=new_df, palette=colors, ax=axes[0])
axes[0].set_title('V11 vs Class Positive Correlation')

sns.boxplot(x="Class", y="V4", data=new_df, palette=colors, ax=axes[1])
axes[1].set_title('V4 vs Class Positive Correlation')

sns.boxplot(x="Class", y="V2", data=new_df, palette=colors, ax=axes[2])
axes[2].set_title('V2 vs Class Positive Correlation')

sns.boxplot(x="Class", y="V19", data=new_df, palette=colors, ax=axes[3])
axes[3].set_title('V19 vs Class Positive Correlation')

plt.show()
```



```
In [37]: from scipy.stats import norm

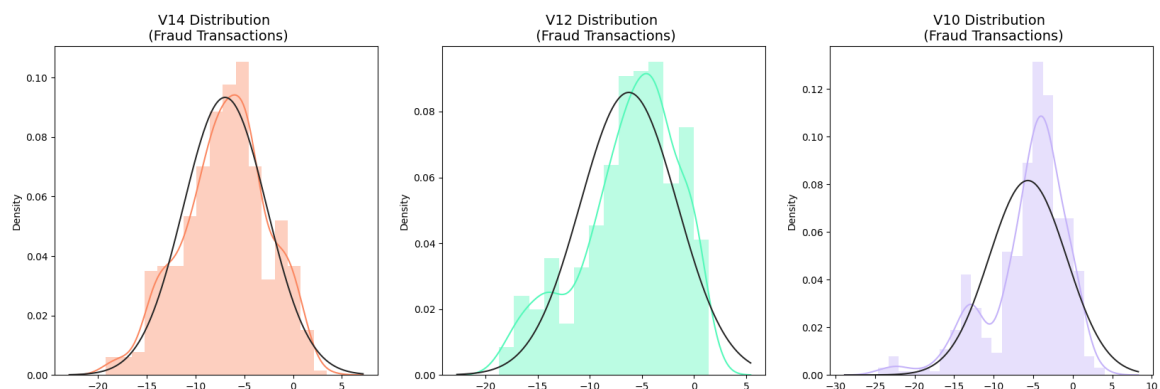
f, (ax1, ax2, ax3) = plt.subplots(1,3, figsize=(20, 6))

v14_fraud_dist = new_df['V14'].loc[new_df['Class'] == 1].values
sns.distplot(v14_fraud_dist,ax=ax1, fit=norm, color='#FB8861')
ax1.set_title('V14 Distribution \n (Fraud Transactions)', fontsize=14)

v12_fraud_dist = new_df['V12'].loc[new_df['Class'] == 1].values
sns.distplot(v12_fraud_dist,ax=ax2, fit=norm, color='#56F9BB')
ax2.set_title('V12 Distribution \n (Fraud Transactions)', fontsize=14)

v10_fraud_dist = new_df['V10'].loc[new_df['Class'] == 1].values
sns.distplot(v10_fraud_dist,ax=ax3, fit=norm, color='#C5B3F9')
ax3.set_title('V10 Distribution \n (Fraud Transactions)', fontsize=14)

plt.show()
```



```
In [38]: f,(ax1, ax2, ax3) = plt.subplots(1, 3, figsize=(20,6))

colors = ['#B3F9C5', '#f9c5b3']
# Boxplots with outliers removed
# Feature V14
sns.boxplot(x="Class", y="V14", data=new_df,ax=ax1, palette=colors)
ax1.set_title("V14 Feature \n Reduction of outliers", fontsize=14)
ax1.annotate('Fewer extreme \n outliers', xy=(0.98, -17.5), xytext=(0, -12),
            arrowprops=dict(facecolor='black'),
            fontsize=14)

# Feature V12
```

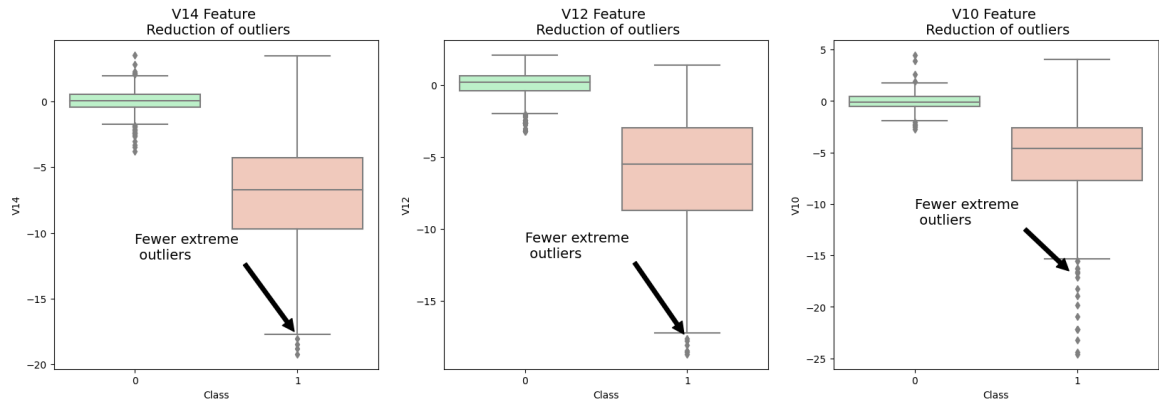
```

sns.boxplot(x="Class", y="V12", data=new_df, ax=ax2, palette=colors)
ax2.set_title("V12 Feature \n Reduction of outliers", fontsize=14)
ax2.annotate('Fewer extreme \n outliers', xy=(0.98, -17.3), xytext=(0, -12),
            arrowprops=dict(facecolor='black'),
            fontsize=14)

# Feature V10
sns.boxplot(x="Class", y="V10", data=new_df, ax=ax3, palette=colors)
ax3.set_title("V10 Feature \n Reduction of outliers", fontsize=14)
ax3.annotate('Fewer extreme \n outliers', xy=(0.95, -16.5), xytext=(0, -12),
            arrowprops=dict(facecolor='black'),
            fontsize=14)

plt.show()

```



```

In [39]: X = new_df.drop('Class', axis=1)
          y = new_df['Class']

# T-SNE Implementation
t0 = time.time()
X_reduced_tsne = TSNE(n_components=2, random_state=42).fit_transform(X.values)
t1 = time.time()
print("T-SNE took {:.2} s".format(t1 - t0))

# PCA Implementation
t0 = time.time()
X_reduced_pca = PCA(n_components=2, random_state=42).fit_transform(X.values)
t1 = time.time()
print("PCA took {:.2} s".format(t1 - t0))

# TruncatedSVD
t0 = time.time()
X_reduced_svd = TruncatedSVD(n_components=2, algorithm='randomized', random_state=42).fit_transform(X.values)
t1 = time.time()
print("Truncated SVD took {:.2} s".format(t1 - t0))

```

T-SNE took 3.2 s

PCA took 0.0 s

Truncated SVD took 0.011 s

```

In [40]: f, (ax1, ax2, ax3) = plt.subplots(1, 3, figsize=(24,6))
          # Labels = ['No Fraud', 'Fraud']
          f.suptitle('Clusters using Dimensionality Reduction', fontsize=14)

          blue_patch = mpatches.Patch(color='#0A0AFF', label='No Fraud')
          red_patch = mpatches.Patch(color='#AF0000', label='Fraud')

          # t-SNE scatter plot
          ax1.scatter(X_reduced_tsne[:,0], X_reduced_tsne[:,1], c=(y == 0), cmap='coolwarm', label='No Fraud')
          ax1.scatter(X_reduced_tsne[:,0], X_reduced_tsne[:,1], c=(y == 1), cmap='coolwarm', label='Fraud',
                    ax1.set_title('t-SNE', fontsize=14)

          ax1.grid(True)

          ax1.legend(handles=[blue_patch, red_patch])

          # PCA scatter plot

```

```

ax2.scatter(X_reduced_pca[:,0], X_reduced_pca[:,1], c=(y == 0), cmap='coolwarm', label='No Fraud',
ax2.scatter(X_reduced_pca[:,0], X_reduced_pca[:,1], c=(y == 1), cmap='coolwarm', label='Fraud', li
ax2.set_title('PCA', fontsize=14)

ax2.grid(True)

ax2.legend(handles=[blue_patch, red_patch])

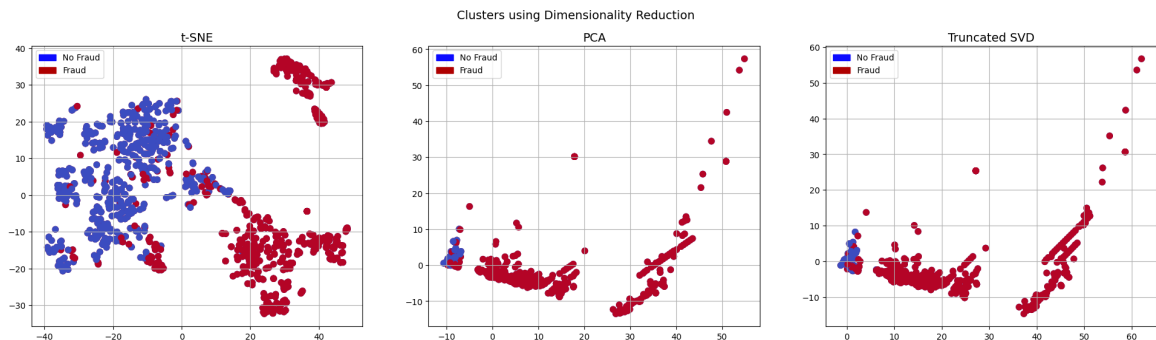
# TruncatedSVD scatter plot
ax3.scatter(X_reduced_svd[:,0], X_reduced_svd[:,1], c=(y == 0), cmap='coolwarm', label='No Fraud',
ax3.scatter(X_reduced_svd[:,0], X_reduced_svd[:,1], c=(y == 1), cmap='coolwarm', label='Fraud', li
ax3.set_title('Truncated SVD', fontsize=14)

ax3.grid(True)

ax3.legend(handles=[blue_patch, red_patch])

plt.show()

```



Classification Models

```

In [41]: # undersampling
X = new_df.drop('Class', axis=1)
y = new_df['Class']

```

```

In [42]: from sklearn.model_selection import train_test_split
# splitting
X_train, X_test, y_train, y_test = train_test_split(X, y, test_size=0.2, random_state=42)
X_train = X_train.values
X_test = X_test.values
y_train = y_train.values
y_test = y_test.values

```

```

In [43]: classifiers = {
    "LogisticRegression": LogisticRegression(),
    "KNearest": KNeighborsClassifier(),
    "Support Vector Classifier": SVC(),
    "DecisionTreeClassifier": DecisionTreeClassifier()
}

```

```

In [44]: from sklearn.model_selection import cross_val_score

for key, classifier in classifiers.items():
    classifier.fit(X_train, y_train)
    training_score = cross_val_score(classifier, X_train, y_train, cv=5)
    print("Classifiers: ", classifier.__class__.__name__, "Has a training score of", round(training_score, 2))

```

```

Classifiers: LogisticRegression Has a training score of 92.0 % accuracy score
Classifiers: KNeighborsClassifier Has a training score of 92.0 % accuracy score
Classifiers: SVC Has a training score of 93.0 % accuracy score
Classifiers: DecisionTreeClassifier Has a training score of 92.0 % accuracy score

```

Cross-Validation

```

In [45]: from sklearn.model_selection import GridSearchCV

# CV logistic regression
log_reg_params = {"penalty": ['l1', 'l2'], 'C': [0.001, 0.01, 0.1, 1, 10, 100, 1000]}
grid_log_reg = GridSearchCV(LogisticRegression(), log_reg_params)

```

```
grid_log_reg.fit(X_train, y_train)

log_reg = grid_log_reg.best_estimator_
log_reg
```

Out[45]: LogisticRegression(C=0.1)

**In a Jupyter environment, please rerun this cell to show the HTML representation or trust the notebook.
On GitHub, the HTML representation is unable to render, please try loading this page with nbviewer.org.**

In [46]:

```
# CV in KNN
kneighbors_params = {"n_neighbors": list(range(2,5,1)), 'algorithm': ['auto', 'ball_tree', 'kd_tree'],
grid_kneighbors = GridSearchCV(KNeighborsClassifier(), kneighbors_params)
grid_kneighbors.fit(X_train, y_train)
kneighbors_neighbors = grid_kneighbors.best_estimator_
```

In [47]:

```
# CV in support vector
svc_params = {'C': [0.5, 0.7, 0.9, 1], 'kernel': ['rbf', 'poly', 'sigmoid', 'linear']}
grid_svc = GridSearchCV(SVC(), svc_params)
grid_svc.fit(X_train, y_train)
svc = grid_svc.best_estimator_
```

In [48]:

```
# CV in DecisionTree

tree_params = {"criterion": ["gini", "entropy"], "max_depth": list(range(2,4,1)),
"min_samples_leaf": list(range(5,7,1))}
grid_tree = GridSearchCV(DecisionTreeClassifier(), tree_params)
grid_tree.fit(X_train, y_train)
tree_clf = grid_tree.best_estimator_
```

In [49]:

```
log_reg_score = cross_val_score(log_reg, X_train, y_train, cv=5)
print('Logistic Regression Cross Validation Score: ', round(log_reg_score.mean() * 100, 2).astype(

kneighbors_score = cross_val_score(kneighbors_neighbors, X_train, y_train, cv=5)
print('KNeares Neighbors Cross Validation Score', round(kneighbors_score.mean() * 100, 2).astype(str) +

svc_score = cross_val_score(svc, X_train, y_train, cv=5)
print('Support Vector Classifier Cross Validation Score', round(svc_score.mean() * 100, 2).astype(

tree_score = cross_val_score(tree_clf, X_train, y_train, cv=5)
print('DecisionTree Classifier Cross Validation Score', round(tree_score.mean() * 100, 2).astype(s
```

Logistic Regression Cross Validation Score: 93.13%
KNeares Neighbors Cross Validation Score 92.62%
Support Vector Classifier Cross Validation Score 92.63%
DecisionTree Classifier Cross Validation Score 92.12%

In [50]:

```
from sklearn.metrics import roc_curve
from sklearn.model_selection import cross_val_predict
# DataFrame with all the scores and the classifiers names.

log_reg_pred = cross_val_predict(log_reg, X_train, y_train, cv=5,
method="decision_function")

kneighbors_pred = cross_val_predict(kneighbors_neighbors, X_train, y_train, cv=5)

svc_pred = cross_val_predict(svc, X_train, y_train, cv=5,
method="decision_function")

tree_pred = cross_val_predict(tree_clf, X_train, y_train, cv=5)
```

In [51]:

```
from sklearn.metrics import roc_auc_score

print('Logistic Regression: ', roc_auc_score(y_train, log_reg_pred))
print('KNeares Neighbors: ', roc_auc_score(y_train, kneares_pred))
print('Support Vector Classifier: ', roc_auc_score(y_train, svc_pred))
print('Decision Tree Classifier: ', roc_auc_score(y_train, tree_pred))
```

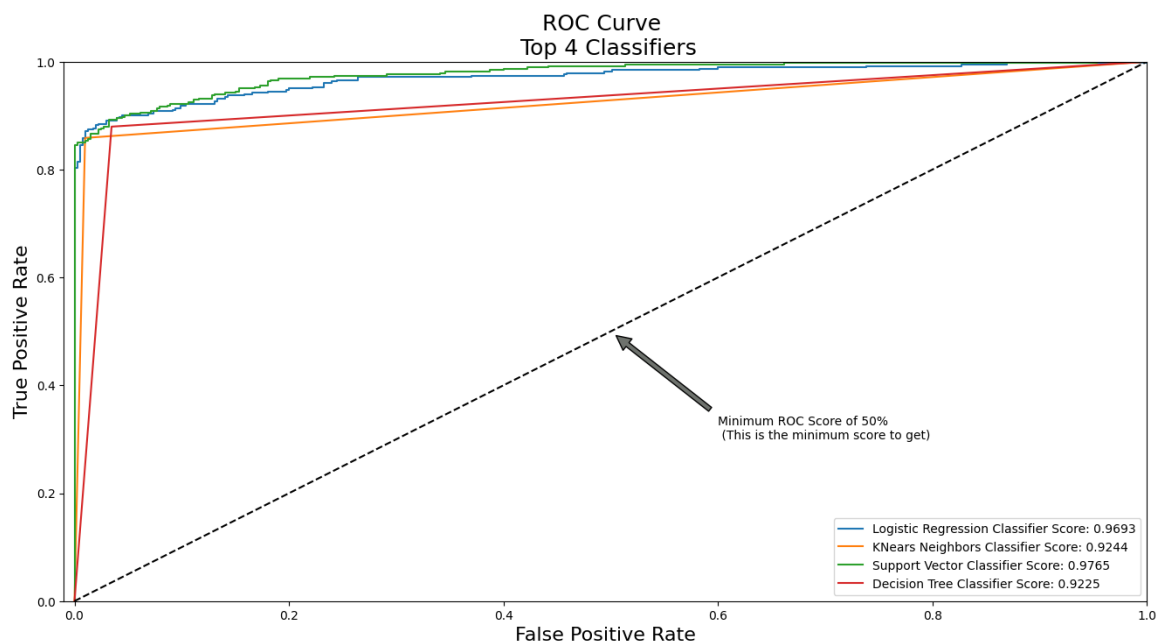
Logistic Regression: 0.9693232499515222
KNeares Neighbors: 0.9243811001228105
Support Vector Classifier: 0.9764591816947837
Decision Tree Classifier: 0.9225066252989464


```
In [52]: # Receiver Operating Characteristic Curve

log_fpr, log_tpr, log_threshold = roc_curve(y_train, log_reg_pred)
knear_fpr, knear_tpr, knear_threshold = roc_curve(y_train, knears_pred)
svc_fpr, svc_tpr, svc_threshold = roc_curve(y_train, svc_pred)
tree_fpr, tree_tpr, tree_threshold = roc_curve(y_train, tree_pred)

def graph_roc_curve_multiple(log_fpr, log_tpr, knear_fpr, knear_tpr, svc_fpr, svc_tpr, tree_fpr, t
plt.figure(figsize=(16,8))
plt.title('ROC Curve \n Top 4 Classifiers', fontsize=18)
plt.plot(log_fpr, log_tpr, label='Logistic Regression Classifier Score: {:.4f}'.format(roc_auc
plt.plot(knear_fpr, knear_tpr, label='KNears Neighbors Classifier Score: {:.4f}'.format(roc_auc
plt.plot(svc_fpr, svc_tpr, label='Support Vector Classifier Score: {:.4f}'.format(roc_auc_sco
plt.plot(tree_fpr, tree_tpr, label='Decision Tree Classifier Score: {:.4f}'.format(roc_auc_sco
plt.plot([0, 1], [0, 1], 'k--')
plt.axis([-0.01, 1, 0, 1])
plt.xlabel('False Positive Rate', fontsize=16)
plt.ylabel('True Positive Rate', fontsize=16)
plt.annotate('Minimum ROC Score of 50% \n (This is the minimum score to get)', xy=(0.5, 0.5),
            arrowprops=dict(facecolor='#6E726D', shrink=0.05),
            )
plt.legend()

graph_roc_curve_multiple(log_fpr, log_tpr, knear_fpr, knear_tpr, svc_fpr, svc_tpr, tree_fpr, tree_
plt.show()
```



Ensemble Methods

```
In [53]: from sklearn.linear_model import LogisticRegression
from sklearn.neighbors import KNeighborsClassifier
from sklearn.ensemble import RandomForestClassifier
from sklearn.model_selection import cross_val_score
```

```
In [54]: clf1 = LogisticRegression()
clf2 = RandomForestClassifier()
clf3 = KNeighborsClassifier()
```

```
In [55]: estimators = [('lr', clf1), ('rf', clf2), ('knn', clf3)]
```

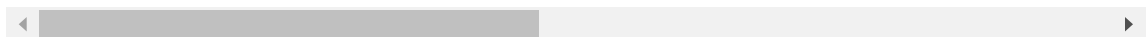
```
In [69]: X = df.iloc[:,0:-1]
y = df.iloc[:, -1]
```

```
In [70]: x
```

```
Out[70]: scaled_amount scaled_time V1 V2 V3 V4 V5 V6
```

242358	0.936212	0.784384	1.968663	-1.209091	-0.511551	-0.983898	-0.970887	0.010964	-1.0830
210127	0.521065	0.624537	1.935025	-0.859136	-1.943729	-0.531092	0.504155	0.706313	-0.1899
195950	-0.043038	0.547469	-3.555417	0.451106	-1.261389	-0.614099	0.674008	1.153825	-4.6222
79470	-0.127157	-0.313150	1.179837	0.362327	0.638986	1.373728	-0.577593	-1.178408	0.1749
100920	-0.279466	-0.200038	1.228461	-1.180218	1.627690	-0.335551	-1.785555	0.901651	-1.7863
...
94600	-0.230699	-0.231958	1.308012	0.322083	-0.045809	0.501397	0.030490	-0.599270	0.0764
51393	1.313491	-0.467216	1.212028	-1.538757	0.685187	-1.393421	-1.739644	0.018674	-1.3282
149398	-0.279746	0.075941	-0.781961	0.092900	-0.423866	-0.329017	2.583220	-0.975481	0.4761
203919	-0.137777	0.591396	-0.212908	1.552188	-0.430994	1.103847	1.468771	-0.822175	1.6886
165116	7.962831	0.382018	-1.587587	-1.957770	0.378966	-1.343979	0.193843	0.033284	-1.7043

284807 rows × 30 columns



In [71]:

y