

In []:

```

1 pip install tensorflow --user
2 !pip install keras
3 !pip install daytime
4 !pip install torch

```

In [1]:

```

1 import pandas as pd
2 import numpy as np
3 import tensorflow as tf
4 import matplotlib.pyplot as plt
5 import seaborn as sns
6 from sklearn.model_selection import train_test_split

```

In [3]:

```

1 from sklearn.preprocessing import StandardScaler
2 from sklearn.metrics import confusion_matrix, recall_score, accuracy_score, precision_score
3 RANDOM_SEED = 2021
4 TEST_PCT = 0.3
5 LABELS = ["Normal", "Fraud"]

```

In [26]:

```

1 dataset = pd.read_csv("E:\Teachning material\Deep learning BE IT 2019 course\creditcard.csv")
2 #dataset.head
3 print(list(dataset.columns))
4 dataset.describe()

```

```

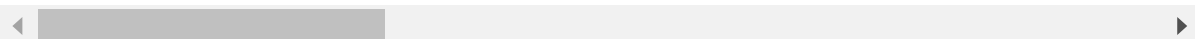
['Time', 'V1', 'V2', 'V3', 'V4', 'V5', 'V6', 'V7', 'V8', 'V9', 'V10', 'V11',
 'V12', 'V13', 'V14', 'V15', 'V16', 'V17', 'V18', 'V19', 'V20', 'V21', 'V22',
 'V23', 'V24', 'V25', 'V26', 'V27', 'V28', 'Amount', 'Class']

```

Out[26]:

	Time	V1	V2	V3	V4	V
count	284807.000000	2.848070e+05	2.848070e+05	2.848070e+05	2.848070e+05	2.848070e+0
mean	94813.859575	3.919560e-15	5.688174e-16	-8.769071e-15	2.782312e-15	-1.552563e-1
std	47488.145955	1.958696e+00	1.651309e+00	1.516255e+00	1.415869e+00	1.380247e+0
min	0.000000	-5.640751e+01	-7.271573e+01	-4.832559e+01	-5.683171e+00	-1.137433e+0
25%	54201.500000	-9.203734e-01	-5.985499e-01	-8.903648e-01	-8.486401e-01	-6.915971e-0
50%	84692.000000	1.810880e-02	6.548556e-02	1.798463e-01	-1.984653e-02	-5.433583e-0
75%	139320.500000	1.315642e+00	8.037239e-01	1.027196e+00	7.433413e-01	6.119264e-0
max	172792.000000	2.454930e+00	2.205773e+01	9.382558e+00	1.687534e+01	3.480167e+0

8 rows × 31 columns



In [6]:

```
1 #check for any nullvalues
2 print("Any nulls in the dataset ",dataset.isnull().values.any() )
3 print('-----')
4 print("No. of unique labels ", len(dataset['Class'].unique()))
5 print("Label values ",dataset.Class.unique())
6 #0 is for normal credit card transaction
7 #1 is for fraudulent credit card transaction
8 print('-----')
9 print("Break down of the Normal and Fraud Transactions")
10 print(pd.value_counts(dataset['Class'], sort = True) )
```

Any nulls in the dataset False

No. of unique labels 2

Label values [0 1]

Break down of the Normal and Fraud Transactions

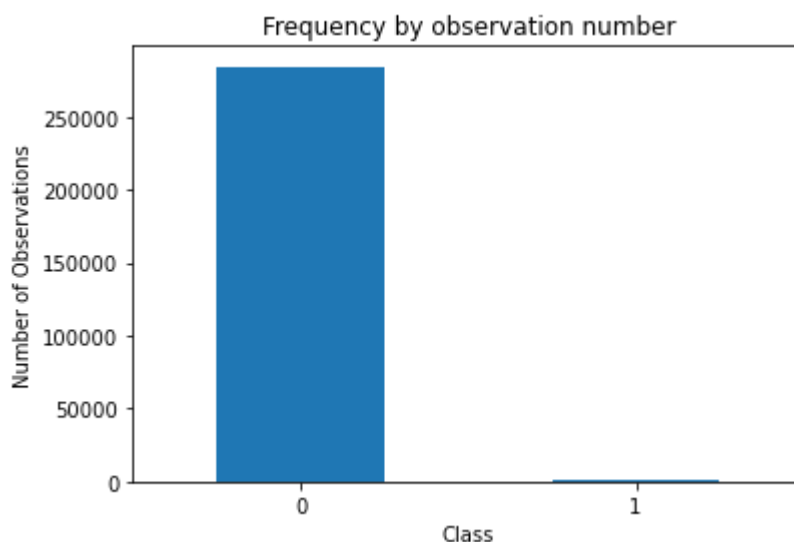
0 284315

1 492

Name: Class, dtype: int64

In [7]:

```
1 #Visualizing the imbalanced dataset
2 count_classes = pd.value_counts(dataset['Class'], sort = True)
3 count_classes.plot(kind = 'bar', rot=0)
4 plt.xticks(range(len(dataset['Class'].unique())), dataset.Class.unique())
5 plt.title("Frequency by observation number")
6 plt.xlabel("Class")
7 plt.ylabel("Number of Observations");
```

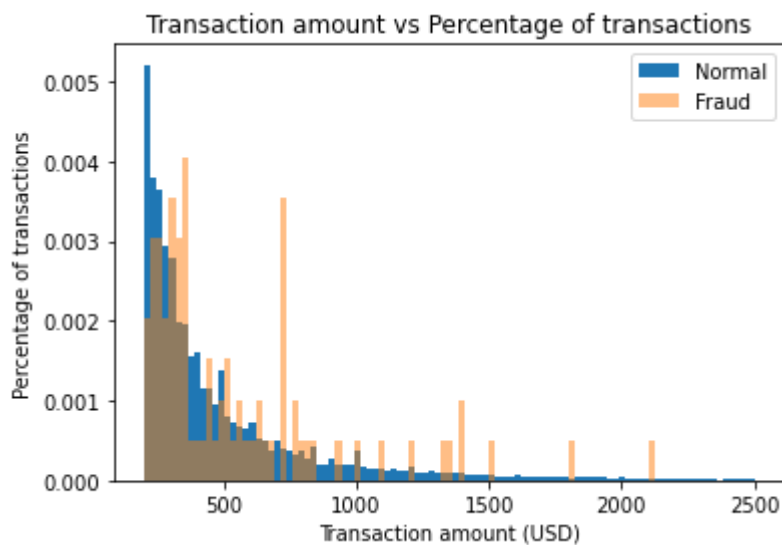


In [8]:

```

1 # Save the normal and fraudulent transactions in separate dataframe
2 normal_dataset = dataset[dataset.Class == 0]
3 fraud_dataset = dataset[dataset.Class == 1]
4 #Visualize transaction amounts for normal and fraudulent transactions
5 bins = np.linspace(200, 2500, 100)
6 plt.hist(normal_dataset.Amount, bins=bins, alpha=1, density=True, label='Normal')
7 plt.hist(fraud_dataset.Amount, bins=bins, alpha=0.5, density=True, label='Fraud')
8 plt.legend(loc='upper right')
9 plt.title("Transaction amount vs Percentage of transactions")
10 plt.xlabel("Transaction amount (USD)")
11 plt.ylabel("Percentage of transactions");
12 plt.show()

```



In []:

```

1 '''Time and Amount are the columns that are not scaled, so applying StandardScaler to c
2 Normalizing the values between 0 and 1 did not work great for the dataset.'''

```

In [9]:

```

1 sc=StandardScaler()
2 dataset['Time'] = sc.fit_transform(dataset['Time'].values.reshape(-1, 1))
3 dataset['Amount'] = sc.fit_transform(dataset['Amount'].values.reshape(-1, 1))

```

In [10]:

```

1  '''The last column in the dataset is our target variable.'''
2
3  raw_data = dataset.values
4  # The last element contains if the transaction is normal which is represented by a 0 or 1
5  labels = raw_data[:, -1]
6  # The other data points are the electrocardiogram data
7  data = raw_data[:, 0:-1]
8  train_data, test_data, train_labels, test_labels = train_test_split(
9      data, labels, test_size=0.2, random_state=2021
10 )

```

In [12]:

```

1  '''Normalize the data to have a value between 0 and 1'''
2
3  min_val = tf.reduce_min(train_data)
4  max_val = tf.reduce_max(train_data)
5  train_data = (train_data - min_val) / (max_val - min_val)
6  test_data = (test_data - min_val) / (max_val - min_val)
7  train_data = tf.cast(train_data, tf.float32)
8  test_data = tf.cast(test_data, tf.float32)

```

In [13]:

```

1  '''Use only normal transactions to train the Autoencoder.
2
3  Normal data has a value of 0 in the target variable. Using the target variable to create
4
5  train_labels = train_labels.astype(bool)
6  test_labels = test_labels.astype(bool)
7
8  #creating normal and fraud datasets
9
10 normal_train_data = train_data[~train_labels]
11 normal_test_data = test_data[~test_labels]
12 fraud_train_data = train_data[train_labels]
13 fraud_test_data = test_data[test_labels]
14 print(" No. of records in Fraud Train Data=",len(fraud_train_data))
15 print(" No. of records in Normal Train data=",len(normal_train_data))
16 print(" No. of records in Fraud Test Data=",len(fraud_test_data))
17 print(" No. of records in Normal Test data=",len(normal_test_data))

```

No. of records in Fraud Train Data= 389
 No. of records in Normal Train data= 227456
 No. of records in Fraud Test Data= 103
 No. of records in Normal Test data= 56859

In [25]:

```

1 nb_epoch = 50
2 batch_size = 64
3 input_dim = normal_train_data.shape[1] #num of columns, 30
4 encoding_dim = 14
5 hidden_dim_1 = int(encoding_dim / 2) #
6 hidden_dim_2=4
7 learning_rate = 1e-7

```

In [27]:

```

1 #input Layer
2 input_layer = tf.keras.layers.Input(shape=(input_dim, ))
3
4 #Encoder
5 encoder = tf.keras.layers.Dense(encoding_dim, activation="tanh",
6                                 activity_regularizer=tf.keras.regularizers.l2(learning_rate))(encoder)
7 encoder=tf.keras.layers.Dropout(0.2)(encoder)
8 encoder = tf.keras.layers.Dense(hidden_dim_1, activation='relu')(encoder)
9 encoder = tf.keras.layers.Dense(hidden_dim_2, activation=tf.nn.leaky_relu)(encoder)
10
11 # Decoder
12 decoder = tf.keras.layers.Dense(hidden_dim_1, activation='relu')(encoder)
13 decoder=tf.keras.layers.Dropout(0.2)(decoder)
14 decoder = tf.keras.layers.Dense(encoding_dim, activation='relu')(decoder)
15 decoder = tf.keras.layers.Dense(input_dim, activation='tanh')(decoder)
16
17 #Autoencoder
18 autoencoder = tf.keras.Model(inputs=input_layer, outputs=decoder)
19 autoencoder.summary()

```

Model: "model"

Layer (type)	Output Shape	Param #
=====		
input_1 (InputLayer)	[(None, 30)]	0
dense (Dense)	(None, 14)	434
dropout (Dropout)	(None, 14)	0
dense_1 (Dense)	(None, 7)	105
dense_2 (Dense)	(None, 4)	32
dense_3 (Dense)	(None, 7)	35
dropout_1 (Dropout)	(None, 7)	0
dense_4 (Dense)	(None, 14)	112
dense_5 (Dense)	(None, 30)	450

```

=====
Total params: 1,168
Trainable params: 1,168
Non-trainable params: 0

```

In [29]:

```
1 """Define the callbacks for checkpoints and early stopping"""
2
3 cp = tf.keras.callbacks.ModelCheckpoint(filepath="autoencoder_fraud.h5",
4                                         mode='min', monitor='val_loss', verbose=2, save_best_only=True)
5 # define our early stopping
6 early_stop = tf.keras.callbacks.EarlyStopping(
7     monitor='val_loss',
8     min_delta=0.0001,
9     patience=10,
10    verbose=1,
11    mode='min',
12    restore_best_weights=True)
```

In [30]:

```
1 #Compile the Autoencoder
2
3 autoencoder.compile(metrics=['accuracy'],
4                        loss='mean_squared_error',
5                        optimizer='adam')
```

In [38]:

```

1 #Train the Autoencoder
2
3 history = autoencoder.fit(normal_train_data, normal_train_data,
4                           epochs=nb_epoch,
5                           batch_size=batch_size,
6                           shuffle=True,
7                           validation_data=(test_data, test_data),
8                           verbose=1,
9                           callbacks=[cp, early_stop]
10                          ).history
11

```

Epoch 1/50

3551/3554 [=====>.] - ETA: 0s - loss: 1.8010e-05 - accuracy: 0.1771

Epoch 1: val_loss did not improve from 0.00002

3554/3554 [=====] - 5s 1ms/step - loss: 1.8009e-05 - accuracy: 0.1770 - val_loss: 5.1378e-05 - val_accuracy: 0.0251

Epoch 2/50

3526/3554 [=====>.] - ETA: 0s - loss: 1.7782e-05 - accuracy: 0.1860

Epoch 2: val_loss did not improve from 0.00002

3554/3554 [=====] - 4s 1ms/step - loss: 1.7783e-05 - accuracy: 0.1863 - val_loss: 3.4659e-05 - val_accuracy: 0.0251

Epoch 3/50

3532/3554 [=====>.] - ETA: 0s - loss: 1.7454e-05 - accuracy: 0.1958

Epoch 3: val_loss did not improve from 0.00002

3554/3554 [=====] - 4s 1ms/step - loss: 1.7450e-05 - accuracy: 0.1958 - val_loss: 3.4017e-05 - val_accuracy: 0.0251

Epoch 4/50

3547/3554 [=====>.] - ETA: 0s - loss: 1.7239e-05 - accuracy: 0.2095

Epoch 4: val_loss did not improve from 0.00002

3554/3554 [=====] - 4s 1ms/step - loss: 1.7239e-05 - accuracy: 0.2095 - val_loss: 3.0229e-05 - val_accuracy: 0.0251

Epoch 5/50

3527/3554 [=====>.] - ETA: 0s - loss: 1.6995e-05 - accuracy: 0.2248

Epoch 5: val_loss did not improve from 0.00002

3554/3554 [=====] - 5s 1ms/step - loss: 1.6998e-05 - accuracy: 0.2248 - val_loss: 2.9758e-05 - val_accuracy: 0.0251

Epoch 6/50

3525/3554 [=====>.] - ETA: 0s - loss: 1.6877e-05 - accuracy: 0.2382

Epoch 6: val_loss did not improve from 0.00002

3554/3554 [=====] - 5s 1ms/step - loss: 1.6878e-05 - accuracy: 0.2382 - val_loss: 2.8776e-05 - val_accuracy: 0.0251

Epoch 7/50

3516/3554 [=====>.] - ETA: 0s - loss: 1.6769e-05 - accuracy: 0.2545

Epoch 7: val_loss did not improve from 0.00002

3554/3554 [=====] - 5s 1ms/step - loss: 1.6757e-05 - accuracy: 0.2545 - val_loss: 2.5803e-05 - val_accuracy: 0.0230

Epoch 8/50

3554/3554 [=====] - ETA: 0s - loss: 1.6662e-05 - accuracy: 0.2617

Epoch 8: val_loss did not improve from 0.00002

```

3554/3554 [=====] - 4s 1ms/step - loss: 1.6662e-05 - accuracy: 0.2617 - val_loss: 2.7398e-05 - val_accuracy: 0.0251
Epoch 9/50
3543/3554 [=====>.] - ETA: 0s - loss: 1.6464e-05 - accuracy: 0.2691
Epoch 9: val_loss did not improve from 0.00002
3554/3554 [=====] - 4s 1ms/step - loss: 1.6460e-05 - accuracy: 0.2690 - val_loss: 2.4035e-05 - val_accuracy: 0.0271
Epoch 10/50
3551/3554 [=====>.] - ETA: 0s - loss: 1.6378e-05 - accuracy: 0.2717
Epoch 10: val_loss did not improve from 0.00002
3554/3554 [=====] - 4s 1ms/step - loss: 1.6380e-05 - accuracy: 0.2717 - val_loss: 2.5148e-05 - val_accuracy: 0.0713
Epoch 11/50
3507/3554 [=====>.] - ETA: 0s - loss: 1.6261e-05 - accuracy: 0.2759
Epoch 11: val_loss did not improve from 0.00002
Restoring model weights from the end of the best epoch: 1.
3554/3554 [=====] - 4s 1ms/step - loss: 1.6261e-05 - accuracy: 0.2757 - val_loss: 2.4422e-05 - val_accuracy: 0.0692
Epoch 11: early stopping

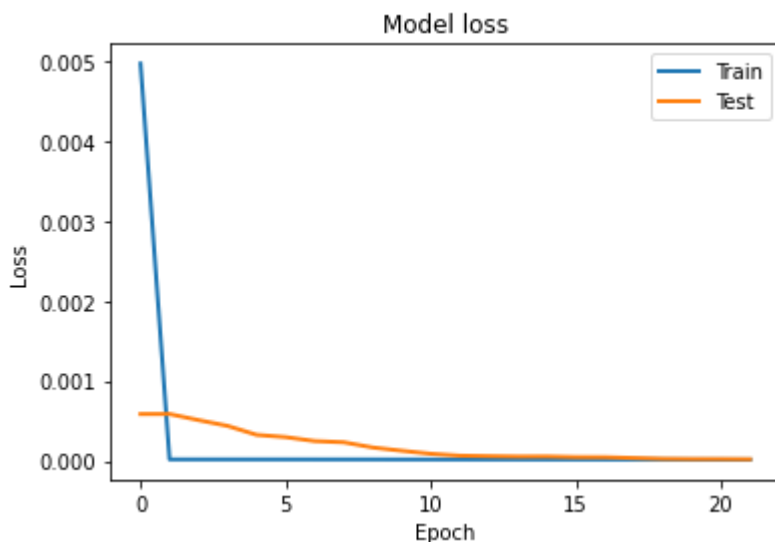
```

In [32]:

```

1 #Plot training and test loss
2
3 plt.plot(history['loss'], linewidth=2, label='Train')
4 plt.plot(history['val_loss'], linewidth=2, label='Test')
5 plt.legend(loc='upper right')
6 plt.title('Model loss')
7 plt.ylabel('Loss')
8 plt.xlabel('Epoch')
9 #plt.ylim(ymin=0.70,ymax=1)
10 plt.show()

```



In [33]:

```

1  """Detect Anomalies on test data
2
3  Anomalies are data points where the reconstruction loss is higher
4
5  To calculate the reconstruction loss on test data,
6  predict the test data and calculate the mean square error between the test data and the
7
8  test_x_predictions = autoencoder.predict(test_data)
9  mse = np.mean(np.power(test_data - test_x_predictions, 2), axis=1)
10 error_df = pd.DataFrame({'Reconstruction_error': mse,
11                          'True_class': test_labels})

```

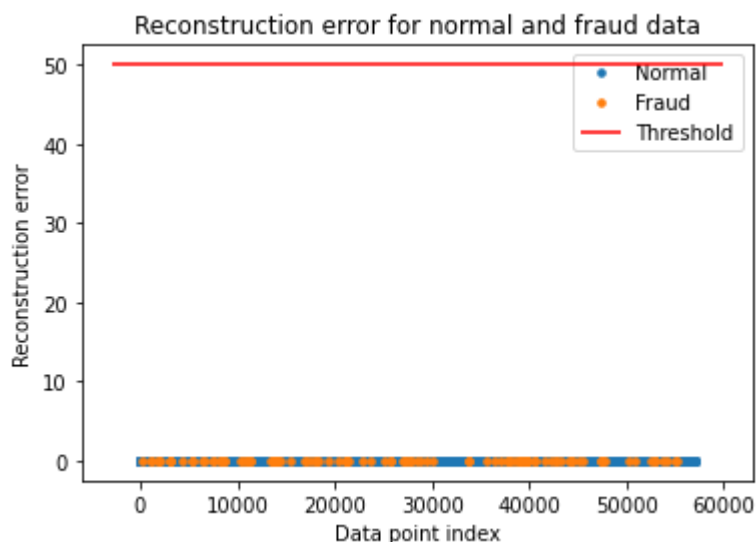
1781/1781 [=====] - 1s 680us/step

In [34]:

```

1  #Plotting the test data points and their respective reconstruction error sets a threshold
2  #if the threshold value needs to be adjusted.
3
4  threshold_fixed = 50
5  groups = error_df.groupby('True_class')
6  fig, ax = plt.subplots()
7  for name, group in groups:
8      ax.plot(group.index, group.Reconstruction_error, marker='o', ms=3.5, linestyle='',
9              label= "Fraud" if name == 1 else "Normal")
10 ax.hlines(threshold_fixed, ax.get_xlim()[0], ax.get_xlim()[1], colors="r", zorder=100,
11 ax.legend()
12 plt.title("Reconstruction error for normal and fraud data")
13 plt.ylabel("Reconstruction error")
14 plt.xlabel("Data point index")
15 plt.show();

```

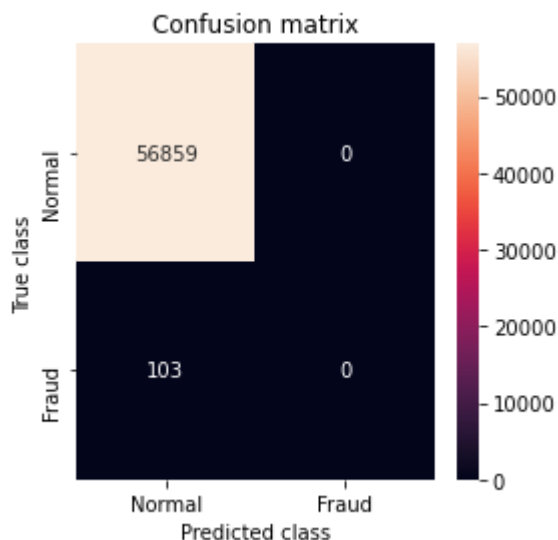


In [35]:

```

1  '''Detect anomalies as points where the reconstruction loss is greater than a fixed thr
2  Here we see that a value of 52 for the threshold will be good.
3
4  Evaluating the performance of the anomaly detection'''
5
6  threshold_fixed =52
7  pred_y = [1 if e > threshold_fixed else 0 for e in error_df.Reconstruction_error.values]
8  error_df['pred'] =pred_y
9  conf_matrix = confusion_matrix(error_df.True_class, pred_y)
10 plt.figure(figsize=(4, 4))
11 sns.heatmap(conf_matrix, xticklabels=LABELS, yticklabels=LABELS, annot=True, fmt="d");
12 plt.title("Confusion matrix")
13 plt.ylabel('True class')
14 plt.xlabel('Predicted class')
15 plt.show()
16 # print Accuracy, precision and recall
17 print(" Accuracy: ",accuracy_score(error_df['True_class'], error_df['pred']))
18 print(" Recall: ",recall_score(error_df['True_class'], error_df['pred']))
19 print(" Precision: ",precision_score(error_df['True_class'], error_df['pred']))

```



Accuracy: 0.9981917769741231

Recall: 0.0

Precision: 0.0

C:\ProgramData\Anaconda3\lib\site-packages\sklearn\metrics_classification.
y:1221: UndefinedMetricWarning: Precision is ill-defined and being set to 0.
0 due to no predicted samples. Use `zero_division` parameter to control this
behavior.

```
_warn_prf(average, modifier, msg_start, len(result))
```

In []:

```
1 '''As our dataset is highly imbalanced, we see a high accuracy but a low recall and pre
2
3 Things to further improve precision and recall would add more relevant features,
4 different architecture for autoencoder, different hyperparameters, or a different algo
```

In []:

```
1 history
```