

Fraud Analytics (CS6890)

Assignment:

4

Title : Fraud Detection Using an Autoencoder and Variational Autoencoder

	Name	Roll Number
	<i>Shreesh Gupta</i>	<i>CS23MTECH12009</i>
	<i>Hrishikesh Hemke</i>	<i>CS23MTECH14003</i>
Team Details :		
	<i>Manan Patel</i>	<i>CS23MTECH14006</i>
	<i>Yug Patel</i>	<i>CS23MTECH14019</i>
	<i>Bhargav Patel</i>	<i>CS23MTECH11026</i>

In [4]:

```
1 import pandas as pd
2 import numpy as np
3 from sklearn.metrics import roc_curve, auc
4 import matplotlib.pyplot as plt
5 import tensorflow as tf
6 import seaborn as sns
7 from sklearn.model_selection import train_test_split
8 from keras.models import Model, load_model
9 from keras.layers import Input, Dense
10 from keras.callbacks import ModelCheckpoint, TensorBoard
11 from keras import regularizers
12 from tensorflow.keras import models, layers
13 from tensorflow.keras.models import Model
14 from sklearn.utils import resample
15 from sklearn.model_selection import train_test_split
16 from sklearn.preprocessing import StandardScaler
17 from tensorflow.keras import layers, models, backend as K
18 from tensorflow.keras.optimizers import Adam
19 from tensorflow.keras.losses import MeanSquaredError
20 from tensorflow.keras.callbacks import EarlyStopping, ReduceLROnPlateau
21 from sklearn.metrics import (roc_curve, auc)
22 from sklearn.manifold import TSNE
```

```
In [5]:  
1 # Load the dataset  
2 data = pd.read_csv('creditcard.csv')  
3  
4 # Display the first few rows of the dataset for a quick overview  
5 data.head()
```

Out[5]:

	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

5 rows × 31 columns

◀ ▶

Data Exploration

Here we are plotting a count plot to visualize the distribution of fraud vs. non-fraud transactions in a dataset.

In [6]:

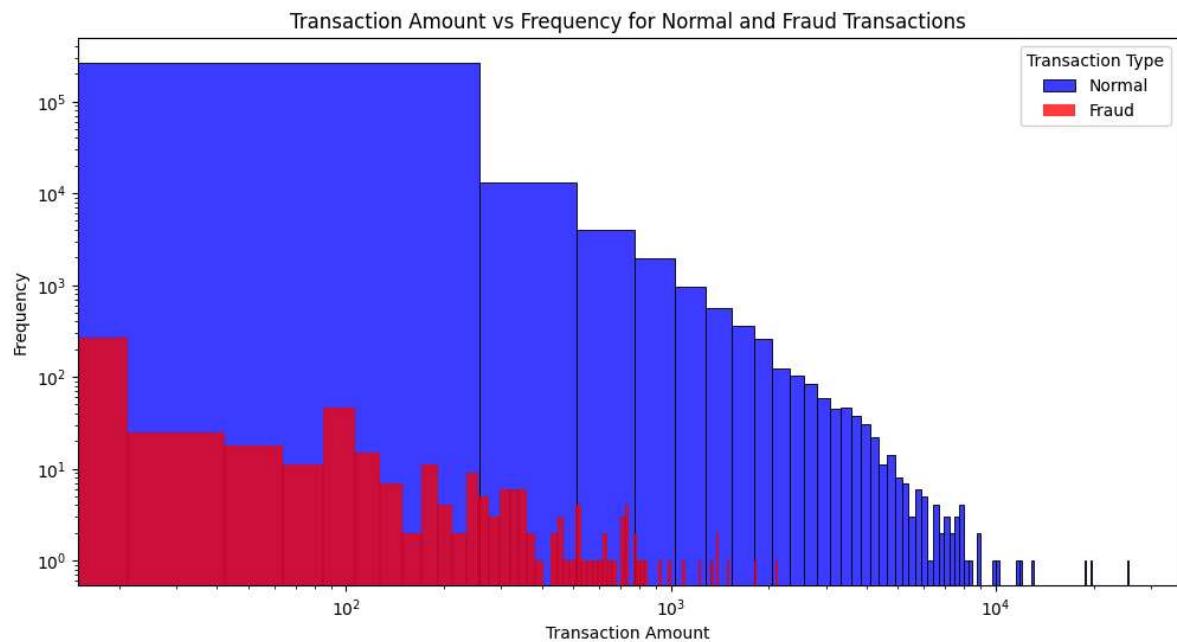
```
1 # Create the count plot
2 fig, ax = plt.subplots(figsize=(9, 6))
3 sns.countplot(x='Class', data=data, ax=ax)
4
5 # Set the title, x-label, and y-label
6 ax.set_title('Count of Fraud vs. Non-Fraud Transactions')
7 ax.set_xlabel('Class (0: Non-Fraud, 1: Fraud)')
8 ax.set_ylabel('Count')
9
10 # Add count labels above the bars
11 for p in ax.patches:
12     ax.annotate(format(p.get_height(), '.0f'),
13                 (p.get_x() + p.get_width() / 2., p.get_height()),
14                 ha = 'center', va = 'center',
15                 xytext = (0, 10),
16                 textcoords = 'offset points')
17
18 # Adjust Layout and display the plot
19 plt.tight_layout()
20 plt.show()
```



Here we are plotting histogram that visually compares the distribution of transaction amounts for normal and fraud transaction. Here we can detect differences between the two classes

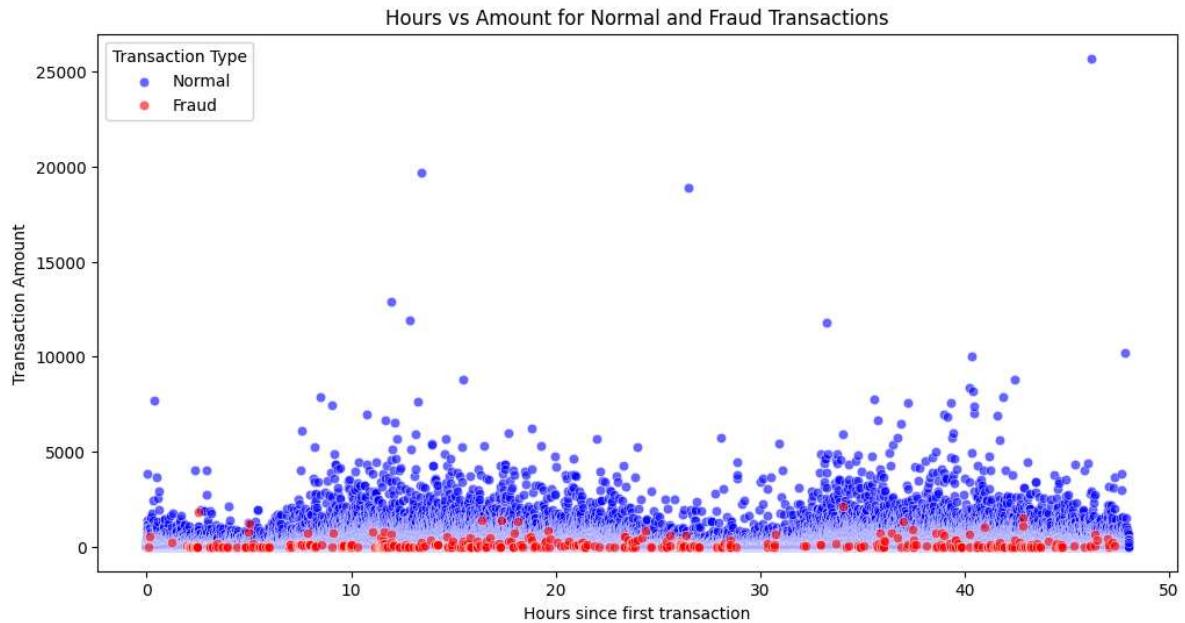
In [7]:

```
1 # Create a figure for plotting
2 plt.figure(figsize=(12, 6))
3
4 # Plotting transaction amount vs frequency for normal transactions
5 sns.histplot(data[data['Class'] == 0]['Amount'], bins=100, color='blue', log=True)
6 sns.histplot(data[data['Class'] == 1]['Amount'], bins=100, color='red', log=True)
7
8 plt.title('Transaction Amount vs Frequency for Normal and Fraud Transactions')
9 plt.xlabel('Transaction Amount')
10 plt.ylabel('Frequency')
11 plt.legend(title='Transaction Type')
12 plt.xscale('log') # Using log scale for better visualization of the spread
13 plt.yscale('log') # Log scale for frequency to handle wide range of values
14 plt.show()
```



In our approach we are converting time into hours and min to do feature selection

```
In [8]:  
1 # Convert 'Time' from seconds to hours to make it more interpretable  
2 data['Hours'] = data['Time'] / 3600  
3  
4 # Create a figure for plotting  
5 plt.figure(figsize=(12, 6))  
6  
7 # Scatter plot of Time vs Amount for normal transactions  
8 sns.scatterplot(x='Hours', y='Amount', data=data[data['Class'] == 0], color='blue')  
9  
10 # Scatter plot of Time vs Amount for fraud transactions  
11 sns.scatterplot(x='Hours', y='Amount', data=data[data['Class'] == 1], color='red')  
12  
13 plt.title('Hours vs Amount for Normal and Fraud Transactions')  
14 plt.xlabel('Hours since first transaction')  
15 plt.ylabel('Transaction Amount')  
16 plt.legend(title='Transaction Type')  
17 plt.show()  
18  
19 data = data.drop('Hours', axis=1)
```



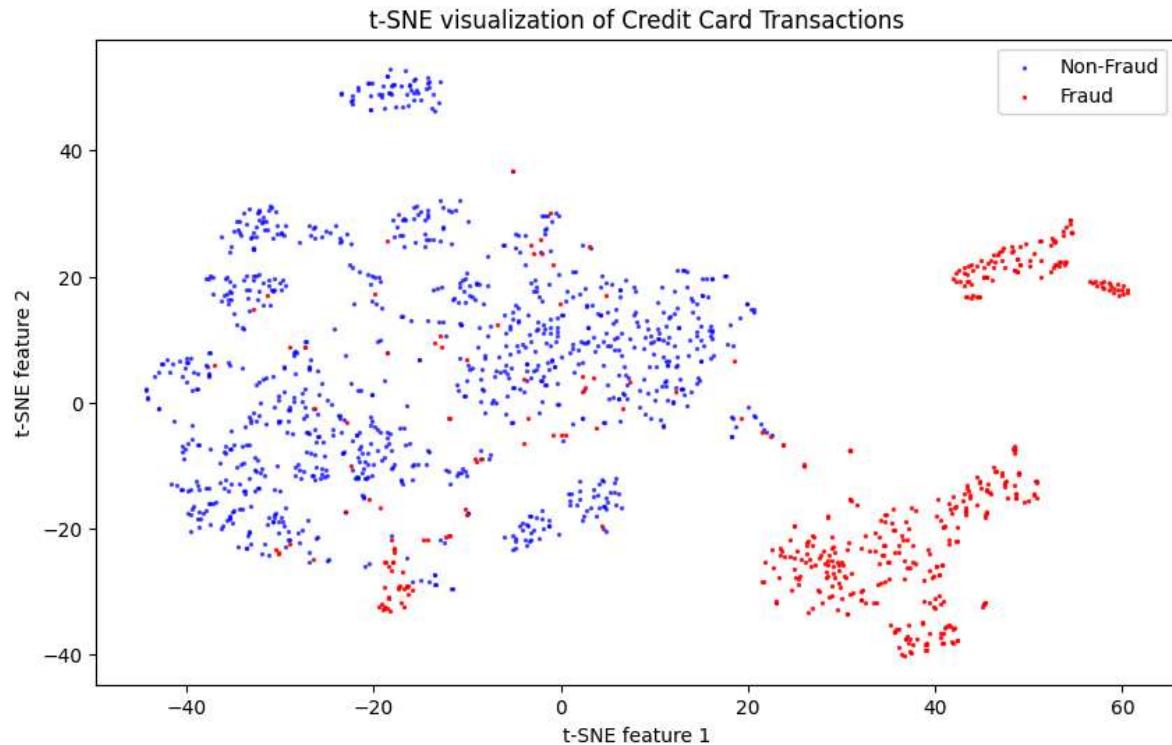
The graph doesn't reveal much about the relationship between fraudulent transactions and the time they occurred. Therefore, we are dropping the 'Time' column.

Data Preprocessing

```
In [ ]: 1 # Convert 'Time' and 'Amount' to log scale to compress dynamic range
2 # Adding 1 before taking the log to avoid Log(0)
3
4 data['Log_Amount'] = np.log(data['Amount'] + 1)
5
6 # Drop the original 'Time' and 'Amount' columns
7 data.drop(['Time', 'Amount'], axis=1, inplace=True)
8
9 # Standardize the newly created log features
10 scaler = StandardScaler()
11 data[['Log_Amount']] = scaler.fit_transform(data[['Log_Amount']])
12
13
14 # Splitting the dataset into training and test sets
15 # Separating features (X) and target variable (y)
16 X = data.drop('Class', axis=1)
17 y = data['Class']
18
19 X_train, X_test, y_train, y_test = train_test_split(X, y, test_size=0.2, r
20
21 # Overview of the split
22 X_train.shape, X_test.shape
```

Out[24]: ((227845, 29), (56962, 29))

```
In [ ]: 1 # Sample for visualizations and model training
2 non_fraud = data[data['Class'] == 0].sample(1000, random_state=42)
3 fraud = data[data['Class'] == 1]
4 df = pd.concat([non_fraud, fraud]).sample(frac=1, random_state=42).reset_index()
5 X = df.drop(['Class'], axis=1).values
6 Y = df['Class'].values
7
8
9 # Assuming tsne_plot is a function that you have defined to plot t-SNE vis
10
11
12 ### Utility Functions
13 def tsne_plot(X, Y, filename):
14     tsne = TSNE(n_components=2, random_state=42)
15     X_tsne = tsne.fit_transform(X)
16     plt.figure(figsize=(10, 6))
17     plt.scatter(X_tsne[Y == 0, 0], X_tsne[Y == 0, 1], label='Non-Fraud', alpha=0.5)
18     plt.scatter(X_tsne[Y == 1, 0], X_tsne[Y == 1, 1], label='Fraud', alpha=0.5)
19     plt.title('t-SNE visualization of Credit Card Transactions')
20     plt.xlabel('t-SNE feature 1')
21     plt.ylabel('t-SNE feature 2')
22     plt.legend()
23     plt.savefig(filename)
24     plt.show()
25
26
27 tsne_plot(X, Y, "original.png")
```



```
In [ ]: 1 # Assuming y_test is a series with binary labels where 1 represents fraud
2 num_fraud_cases = sum(y_test)
3 num_non_fraud_cases = len(y_test) - num_fraud_cases
4
5 # Filter the test set to include only the fraud cases
6 X_test_fraud = X_test[y_test == 1]
7
8 # Sample an equal number of non-fraud cases from the test set
9 X_test_non_fraud = X_test[y_test == 0].sample(n=num_fraud_cases, random_state=42)
10
11 # Combine the fraud and non-fraud cases to form the final test set
12 X_test_final = pd.concat([X_test_fraud, X_test_non_fraud])
13 y_test_final = pd.Series(np.where(X_test_final.index.isin(X_test_fraud.index), 1, 0))
14
15 # Verify the class distribution in the final test set
16 print("Class distribution in the final test set:")
17 print(y_test_final.value_counts())
18
19 # Optionally, shuffle the test set
20 # To shuffle the data correctly, use the sklearn shuffle utility
21 from sklearn.utils import shuffle
22 X_test_final, y_test_final = shuffle(X_test_final, y_test_final, random_state=42)
```

Class distribution in the final test set:

```
1    98
0    98
Name: count, dtype: int64
```

AutoEncoder

Architecher

```
In [ ]: 1 input_dim = X_train.shape[1]
2 encoding_dim = 14
3 input_layer = Input(shape=(input_dim, ))
4 encoder = Dense(encoding_dim, activation="tanh",
5                  activity_regularizer=regularizers.l1(10e-5))(input_layer)
6 encoder = Dense(int(encoding_dim / 2), activation="relu")(encoder)
7 decoder = Dense(int(encoding_dim / 2), activation='tanh')(encoder)
8 decoder = Dense(input_dim, activation='relu')(decoder)
9 autoencoder = Model(inputs=input_layer, outputs=decoder)
```

```
In [ ]: 1 autoencoder.summary()
```

Model: "functional_3"

Layer (type)	Output Shape	Param #
input_layer_1 (InputLayer)	(None, 29)	0
dense_4 (Dense)	(None, 14)	420
dense_5 (Dense)	(None, 7)	105
dense_6 (Dense)	(None, 7)	56
dense_7 (Dense)	(None, 29)	232

Total params: 813 (3.18 KB)

Trainable params: 813 (3.18 KB)

Non-trainable params: 0 (0.00 B)

Training

In []:

```
1 nb_epoch = 50
2 batch_size = 256
3 autoencoder.compile(optimizer='adam',
4                      loss='mean_squared_error',
5                      metrics=['accuracy'])
6 checkpointer = ModelCheckpoint(filepath="model.keras",
7                                 verbose=0,
8                                 save_best_only=True)
9 tensorboard = TensorBoard(log_dir='./logs',
10                           histogram_freq=0,
11                           write_graph=True,
12                           write_images=True)
13 history = autoencoder.fit(X_train, X_train,
14                           epochs=nb_epoch,
15                           batch_size=batch_size,
16                           shuffle=True,
17                           validation_data=(X_test, X_test),
18                           verbose=1,
19                           callbacks=[checkpointer, tensorboard]).history
```

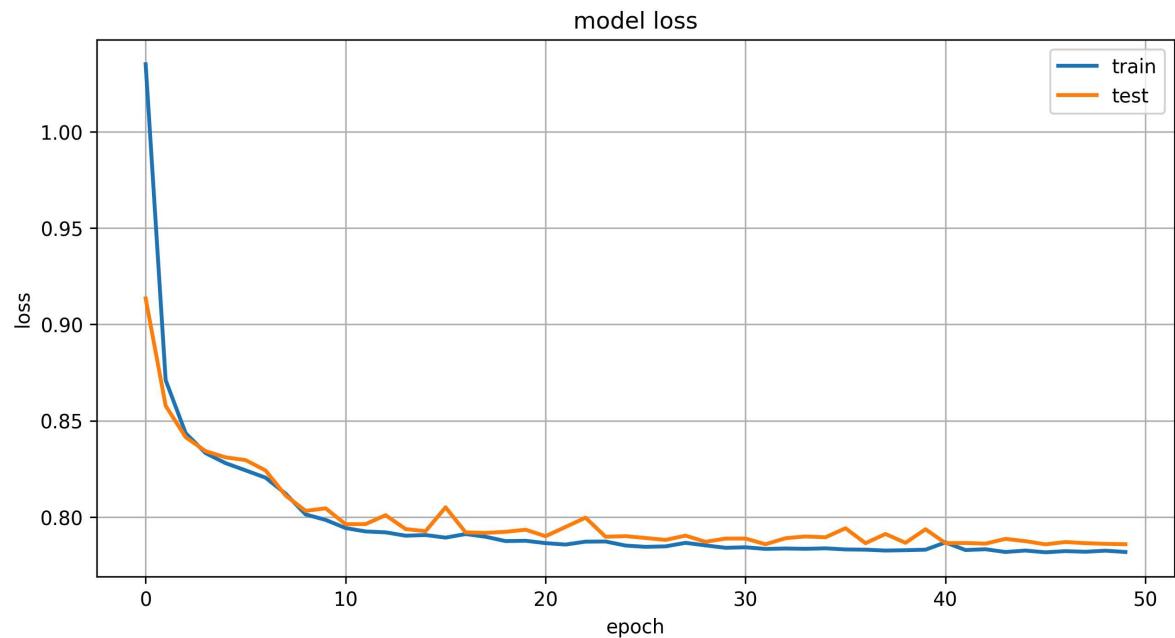
```
Epoch 1/50
891/891 2s 1ms/step - accuracy: 0.2410 - loss: 1.1436 -
val_accuracy: 0.5526 - val_loss: 0.9135
Epoch 2/50
891/891 1s 1ms/step - accuracy: 0.5678 - loss: 0.8721 -
val_accuracy: 0.5784 - val_loss: 0.8578
Epoch 3/50
891/891 1s 1ms/step - accuracy: 0.5857 - loss: 0.8364 -
val_accuracy: 0.5973 - val_loss: 0.8414
Epoch 4/50
891/891 1s 1ms/step - accuracy: 0.6003 - loss: 0.8548 -
val_accuracy: 0.6025 - val_loss: 0.8342
Epoch 5/50
891/891 1s 1ms/step - accuracy: 0.6031 - loss: 0.8376 -
val_accuracy: 0.6059 - val_loss: 0.8310
Epoch 6/50
891/891 1s 1ms/step - accuracy: 0.6024 - loss: 0.8271 -
val_accuracy: 0.5960 - val_loss: 0.8296
Epoch 7/50
891/891 1s 1ms/step - accuracy: 0.6027 - loss: 0.8299 -
val_accuracy: 0.5986 - val_loss: 0.8241
Epoch 8/50
891/891 1s 1ms/step - accuracy: 0.6070 - loss: 0.8171 -
val_accuracy: 0.6212 - val_loss: 0.8111
Epoch 9/50
891/891 1s 1ms/step - accuracy: 0.6268 - loss: 0.8324 -
val_accuracy: 0.6299 - val_loss: 0.8032
Epoch 10/50
891/891 1s 1ms/step - accuracy: 0.6305 - loss: 0.7893 -
val_accuracy: 0.6369 - val_loss: 0.8045
Epoch 11/50
891/891 1s 1ms/step - accuracy: 0.6327 - loss: 0.7905 -
val_accuracy: 0.6336 - val_loss: 0.7963
Epoch 12/50
891/891 1s 1ms/step - accuracy: 0.6320 - loss: 0.8069 -
val_accuracy: 0.6317 - val_loss: 0.7963
Epoch 13/50
891/891 1s 1ms/step - accuracy: 0.6319 - loss: 0.7992 -
val_accuracy: 0.6247 - val_loss: 0.8010
Epoch 14/50
891/891 1s 1ms/step - accuracy: 0.6318 - loss: 0.8007 -
val_accuracy: 0.6332 - val_loss: 0.7937
Epoch 15/50
891/891 1s 1ms/step - accuracy: 0.6331 - loss: 0.7951 -
val_accuracy: 0.6276 - val_loss: 0.7926
Epoch 16/50
891/891 1s 1ms/step - accuracy: 0.6332 - loss: 0.7868 -
val_accuracy: 0.6008 - val_loss: 0.8050
Epoch 17/50
891/891 1s 1ms/step - accuracy: 0.6303 - loss: 0.7886 -
val_accuracy: 0.6367 - val_loss: 0.7920
Epoch 18/50
891/891 1s 1ms/step - accuracy: 0.6328 - loss: 0.7957 -
val_accuracy: 0.6202 - val_loss: 0.7917
Epoch 19/50
891/891 1s 1ms/step - accuracy: 0.6325 - loss: 0.7770 -
val_accuracy: 0.6280 - val_loss: 0.7923
```

```
Epoch 20/50
891/891 1s 1ms/step - accuracy: 0.6325 - loss: 0.7985 -
val_accuracy: 0.6348 - val_loss: 0.7934
Epoch 21/50
891/891 1s 1ms/step - accuracy: 0.6312 - loss: 0.7898 -
val_accuracy: 0.6307 - val_loss: 0.7900
Epoch 22/50
891/891 1s 1ms/step - accuracy: 0.6354 - loss: 0.7926 -
val_accuracy: 0.6312 - val_loss: 0.7948
Epoch 23/50
891/891 1s 1ms/step - accuracy: 0.6351 - loss: 0.7759 -
val_accuracy: 0.6207 - val_loss: 0.7997
Epoch 24/50
891/891 1s 1ms/step - accuracy: 0.6291 - loss: 0.7906 -
val_accuracy: 0.6415 - val_loss: 0.7897
Epoch 25/50
891/891 2s 2ms/step - accuracy: 0.6355 - loss: 0.7848 -
val_accuracy: 0.6318 - val_loss: 0.7901
Epoch 26/50
891/891 1s 1ms/step - accuracy: 0.6345 - loss: 0.7805 -
val_accuracy: 0.6361 - val_loss: 0.7891
Epoch 27/50
891/891 1s 1ms/step - accuracy: 0.6337 - loss: 0.7805 -
val_accuracy: 0.6282 - val_loss: 0.7881
Epoch 28/50
891/891 1s 1ms/step - accuracy: 0.6331 - loss: 0.7641 -
val_accuracy: 0.6377 - val_loss: 0.7903
Epoch 29/50
891/891 1s 1ms/step - accuracy: 0.6349 - loss: 0.8123 -
val_accuracy: 0.6427 - val_loss: 0.7871
Epoch 30/50
891/891 1s 1ms/step - accuracy: 0.6350 - loss: 0.7750 -
val_accuracy: 0.6381 - val_loss: 0.7888
Epoch 31/50
891/891 1s 1ms/step - accuracy: 0.6331 - loss: 0.7925 -
val_accuracy: 0.6335 - val_loss: 0.7888
Epoch 32/50
891/891 1s 1ms/step - accuracy: 0.6368 - loss: 0.7784 -
val_accuracy: 0.6353 - val_loss: 0.7859
Epoch 33/50
891/891 1s 1ms/step - accuracy: 0.6359 - loss: 0.7775 -
val_accuracy: 0.6295 - val_loss: 0.7890
Epoch 34/50
891/891 1s 1ms/step - accuracy: 0.6343 - loss: 0.7783 -
val_accuracy: 0.6331 - val_loss: 0.7899
Epoch 35/50
891/891 1s 1ms/step - accuracy: 0.6347 - loss: 0.7829 -
val_accuracy: 0.6325 - val_loss: 0.7895
Epoch 36/50
891/891 1s 1ms/step - accuracy: 0.6360 - loss: 0.7927 -
val_accuracy: 0.6208 - val_loss: 0.7942
Epoch 37/50
891/891 1s 1ms/step - accuracy: 0.6363 - loss: 0.7878 -
val_accuracy: 0.6329 - val_loss: 0.7864
Epoch 38/50
891/891 1s 1ms/step - accuracy: 0.6347 - loss: 0.7755 -
val_accuracy: 0.6378 - val_loss: 0.7912
```

```
Epoch 39/50
891/891 1s 1ms/step - accuracy: 0.6367 - loss: 0.7818 -
val_accuracy: 0.6391 - val_loss: 0.7866
Epoch 40/50
891/891 1s 1ms/step - accuracy: 0.6362 - loss: 0.7831 -
val_accuracy: 0.6282 - val_loss: 0.7936
Epoch 41/50
891/891 1s 1ms/step - accuracy: 0.6302 - loss: 0.8044 -
val_accuracy: 0.6380 - val_loss: 0.7865
Epoch 42/50
891/891 1s 1ms/step - accuracy: 0.6369 - loss: 0.7820 -
val_accuracy: 0.6310 - val_loss: 0.7865
Epoch 43/50
891/891 1s 1ms/step - accuracy: 0.6323 - loss: 0.7914 -
val_accuracy: 0.6397 - val_loss: 0.7862
Epoch 44/50
891/891 1s 1ms/step - accuracy: 0.6364 - loss: 0.7846 -
val_accuracy: 0.6289 - val_loss: 0.7886
Epoch 45/50
891/891 1s 1ms/step - accuracy: 0.6392 - loss: 0.7938 -
val_accuracy: 0.6341 - val_loss: 0.7875
Epoch 46/50
891/891 1s 1ms/step - accuracy: 0.6354 - loss: 0.7919 -
val_accuracy: 0.6407 - val_loss: 0.7858
Epoch 47/50
891/891 1s 1ms/step - accuracy: 0.6375 - loss: 0.7925 -
val_accuracy: 0.6249 - val_loss: 0.7870
Epoch 48/50
891/891 1s 1ms/step - accuracy: 0.6367 - loss: 0.7788 -
val_accuracy: 0.6507 - val_loss: 0.7864
Epoch 49/50
891/891 1s 1ms/step - accuracy: 0.6376 - loss: 0.7799 -
val_accuracy: 0.6307 - val_loss: 0.7861
Epoch 50/50
891/891 1s 1ms/step - accuracy: 0.6383 - loss: 0.7836 -
val_accuracy: 0.6310 - val_loss: 0.7859
```

In []: 1 autoencoder = load_model('model.keras')

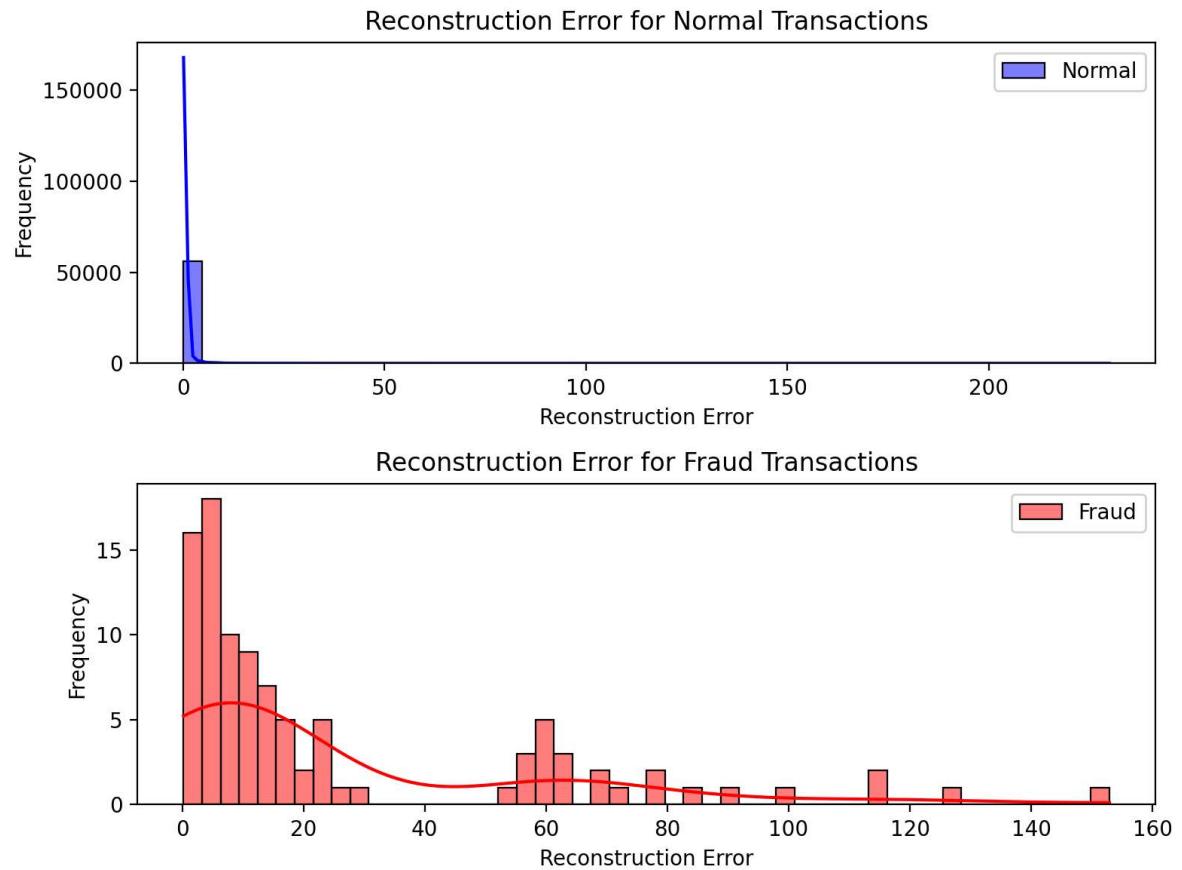
```
In [ ]: 1 plt.figure(figsize=(10, 5), dpi=300) # Setting the figure size and dpi fo  
2 plt.plot(history['loss'], label='train', linewidth=2)  
3 plt.plot(history['val_loss'], label='test', linewidth=2)  
4 plt.title('model loss')  
5 plt.ylabel('loss')  
6 plt.xlabel('epoch')  
7 plt.legend(loc='upper right')  
8 plt.grid(True)  
9 plt.show()
```



```
In [ ]: 1 predictions = autoencoder.predict(X_test)  
2 mse = np.mean(np.power(X_test - predictions, 2), axis=1)
```

1781/1781 ━━━━━━ 1s 473us/step

```
In [ ]: 1 # Separate the reconstruction error by class
2 mse_normal = mse[y_test == 0]
3 mse_fraud = mse[y_test == 1]
4
5 # Using the same mock data for demonstration purposes
6 # Normally you would replace 'mse_normal' and 'mse_fraud' with your actual
7
8 # Create subplots
9 fig, ax = plt.subplots(2, 1, figsize=(8, 6), dpi=200)
10
11 # Histogram of reconstruction errors for normal transactions
12 sns.histplot(mse_normal, bins=50, kde=True, color='blue', label='Normal',
13 ax[0].set_title('Reconstruction Error for Normal Transactions')
14 ax[0].set_xlabel('Reconstruction Error')
15 ax[0].set_ylabel('Frequency')
16 ax[0].legend()
17
18 # Histogram of reconstruction errors for fraud transactions
19 sns.histplot(mse_fraud, bins=50, kde=True, color='red', label='Fraud', ax=
20 ax[1].set_title('Reconstruction Error for Fraud Transactions')
21 ax[1].set_xlabel('Reconstruction Error')
22 ax[1].set_ylabel('Frequency')
23 ax[1].legend()
24
25 plt.tight_layout()
26 plt.show()
```

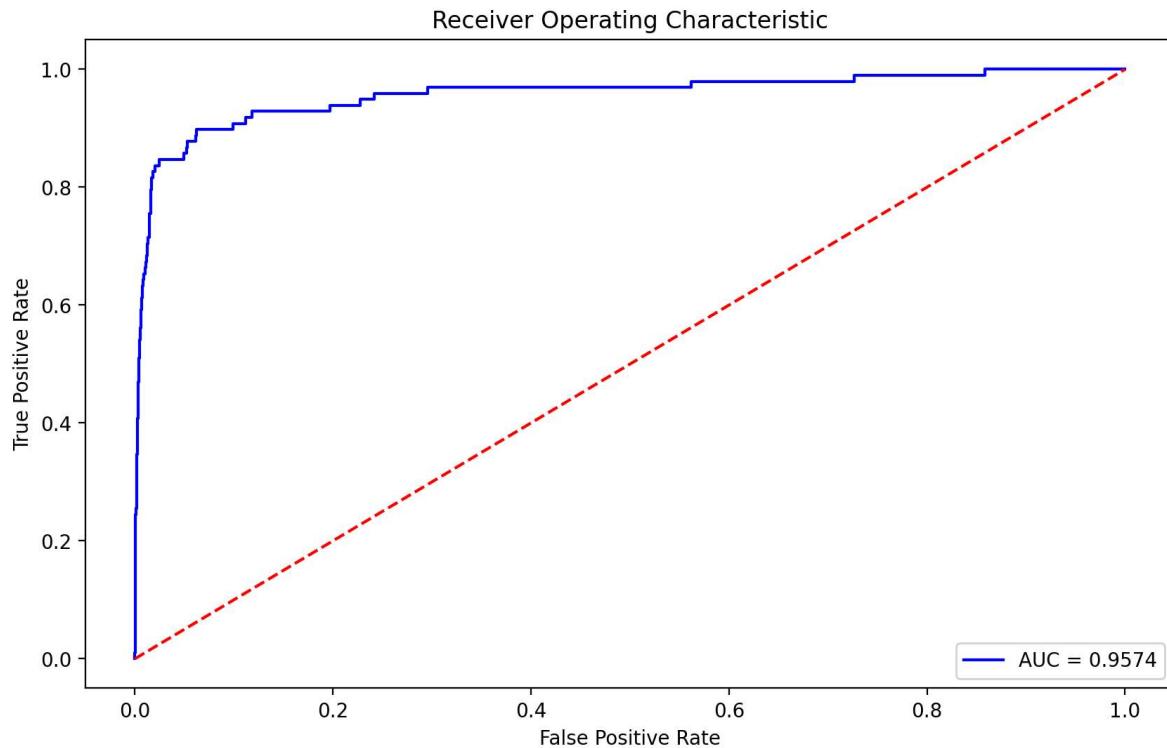


Results

ROC curve

The ROC curve below represents the performance of a classification model, with an AUC (Area Under the Curve) of 0.9574, indicating a high level of accuracy. The true positive rate is high across different thresholds, while maintaining a low false positive rate. This model is effective at distinguishing between the classes (e.g., fraud and non-fraud).

```
In [ ]: 1 # Compute ROC curve and ROC area for the test set
2 fpr, tpr, thresholds = roc_curve(y_test, mse)
3 roc_auc = auc(fpr, tpr)
4
5 plt.figure(figsize=(10, 6), dpi=200)
6 plt.plot(fpr, tpr, color='blue', label=f'AUC = {roc_auc:.4f}')
7 plt.plot([0, 1], [0, 1], color='red', linestyle='--')
8 plt.title('Receiver Operating Characteristic')
9 plt.xlabel('False Positive Rate')
10 plt.ylabel('True Positive Rate')
11 plt.legend(loc='lower right')
12 plt.show()
```



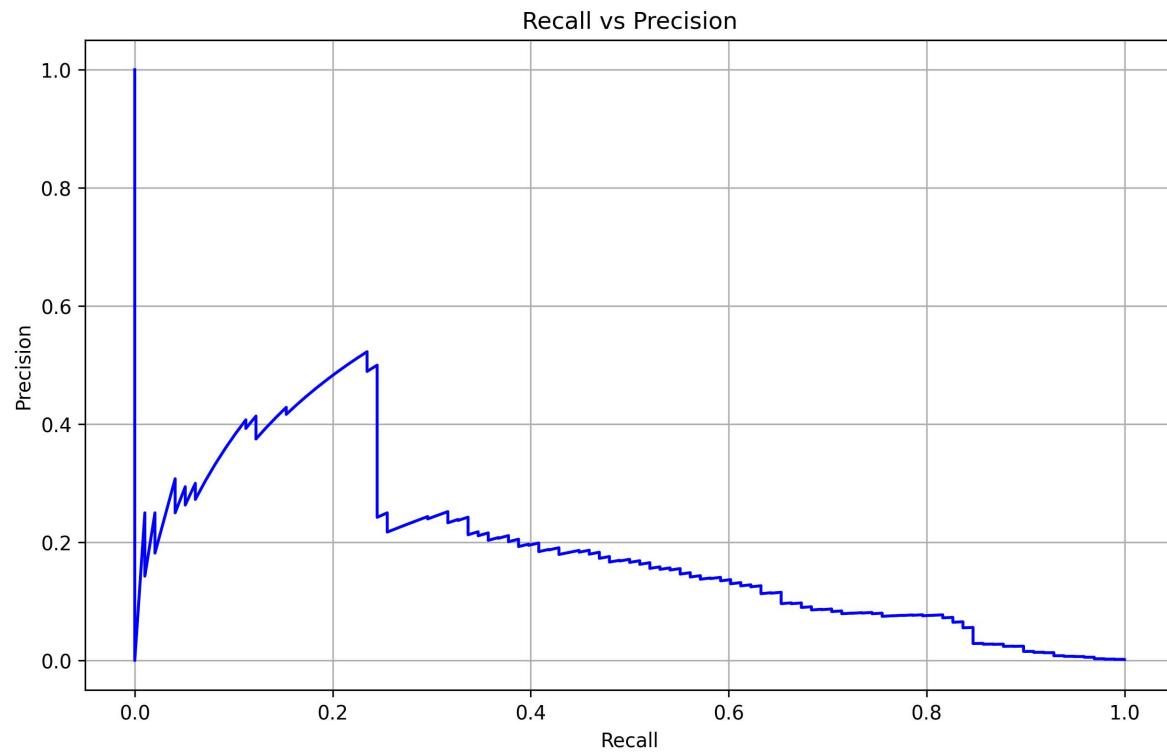
```
In [ ]: 1 from sklearn.metrics import (confusion_matrix, precision_recall_curve, au
2                                         roc_curve, recall_score, classification_report,
3                                         precision_recall_fscore_support)
```

Recall vs Precision

The Precision-Recall curve below shows the trade-off between precision and recall for a classification model at different thresholds. High precision near the y-axis suggests that the model is accurate when it predicts positive classes, but recall drops quickly, indicating it doesn't capture all positive cases. The model is precise but not highly sensitive to all positives.

In []:

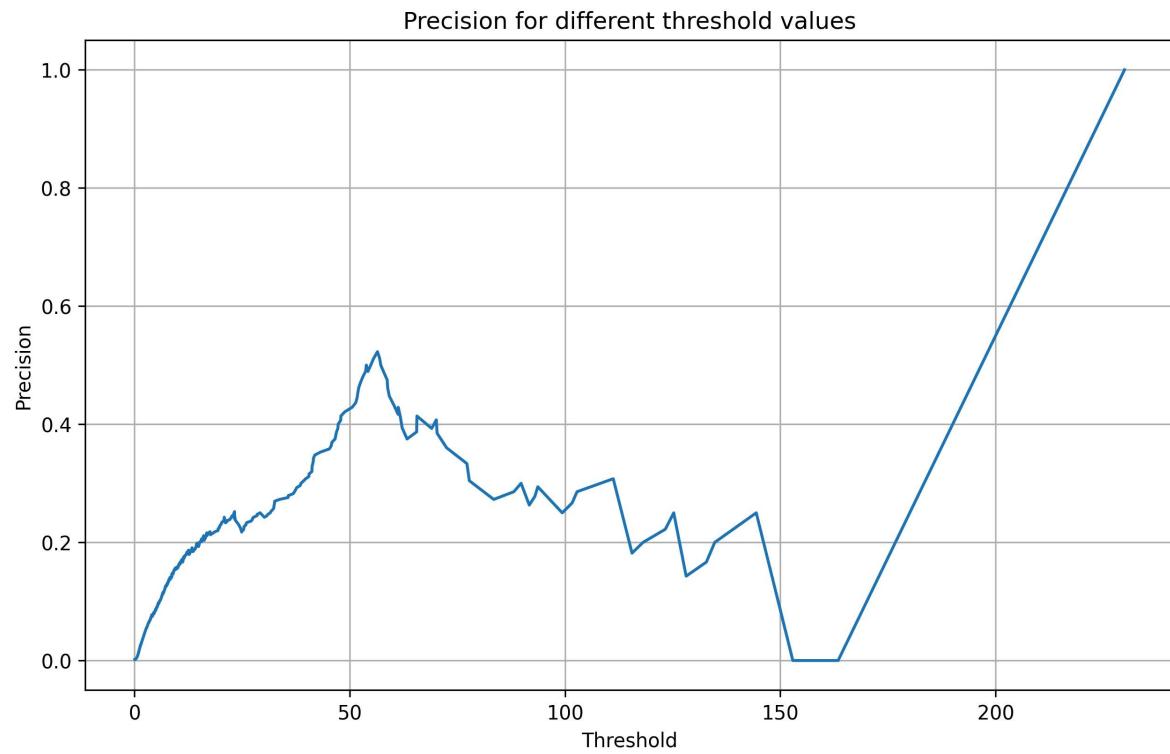
```
1 precision, recall, th = precision_recall_curve(y_test, mse)
2 plt.figure(figsize=(10, 6), dpi=300)
3 plt.plot(recall, precision, 'b')
4 plt.title('Recall vs Precision')
5 plt.xlabel('Recall')
6 plt.ylabel('Precision')
7 plt.grid(True)
8 plt.show()
```



Precision for different threshold values

The graph shows how the precision of a predictive model varies with different threshold values. As the threshold increases, precision fluctuates and tends to decrease before shooting up to perfect precision at the highest threshold. This suggests that the model becomes very selective, possibly only predicting positives when very sure, which may not be practical if high recall is also desired.

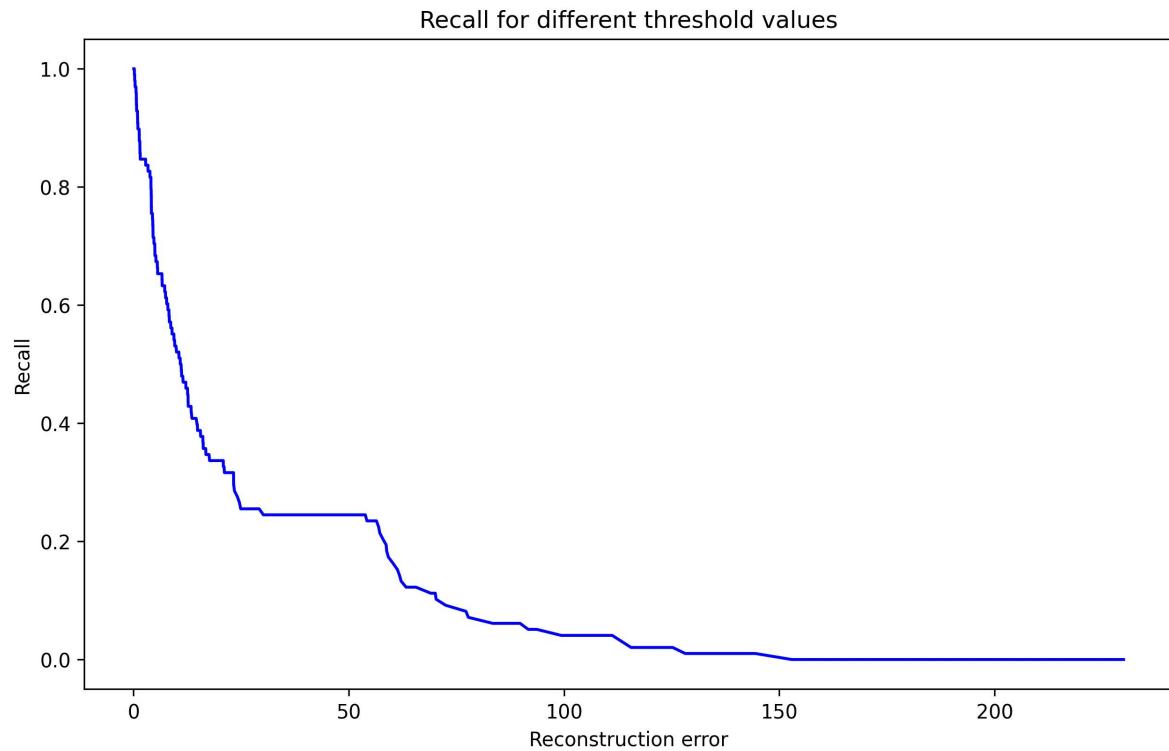
```
In [ ]: 1 plt.figure(figsize=(10, 6), dpi=300)
2 plt.plot(th, precision[1:])
3 plt.title('Precision for different threshold values')
4 plt.xlabel('Threshold')
5 plt.ylabel('Precision')
6 plt.grid(True)
7 plt.show()
```



Threshold-Recall curve

The graph illustrates recall as a function of varying reconstruction error thresholds in a predictive model. Initially, recall is high, indicating most positive cases are captured, but it declines sharply as the threshold increases, implying fewer positives are detected. At high threshold levels, recall is low, suggesting many positive instances are missed by the model.

```
In [ ]: 1 plt.figure(figsize=(10, 6), dpi=300)
2 plt.plot(th, recall[1:], 'b', label='Threshold-Recall curve')
3 plt.title('Recall for different threshold values')
4 plt.xlabel('Reconstruction error')
5 plt.ylabel('Recall')
6 plt.show()
```



F1 SCORE

The calculation of the optimal threshold that maximizes the F1 score involves these concise mathematical steps:

1. Precision and Recall Calculation:

- Generate a set of predictions or error scores (S).
- Use the function `precision_recall_curve(y, S)` to obtain precision ((P)) and recall ((R)) values at various thresholds ((T)) derived from (S).

2. F1 Score Calculation for Each Threshold:

- For each threshold (t) in (T), calculate the F1 score:
$$F1(t) = \frac{2 \times P(t) \times R(t)}{P(t) + R(t)}$$
- This formula ensures that the F1 score is computed for each possible threshold based on the corresponding precision and recall values.

3. Optimal Threshold Determination:

- Identify the threshold (t_{opt}) that maximizes ($F1(t)$):
$$t_{\text{opt}} = \arg\max_t F1(t)$$
- Here, ($\arg\max$) is the operation that finds the argument (threshold) that results in the maximum F1 score.

4. Using the Optimal Threshold:

- With t_{opt} known, classify future data points by comparing their score (s) against t_{opt} .

This process effectively balances precision and recall, yielding the most reliable threshold for making binary decisions in various applications, particularly when handling skewed datasets.

In []:

```

1
2
3 predictions = autoencoder.predict(X_test_final)
4
5 mse_final = np.mean(np.power(X_test_final - predictions, 2), axis=1)
6
7
8
9 precision, recall, thresholds = precision_recall_curve(y_test_final, mse_f
10
11
12 # Calculate F1 scores for each possible threshold
13 f1_scores = 2 * (precision * recall) / (precision + recall)
14
15 # Find the index of the maximum F1 score
16 optimal_idx = np.nanargmax(f1_scores)
17
18
19 optimal_threshold = thresholds[optimal_idx]
20 optimal_f1 = f1_scores[optimal_idx]
21
22 print("Optimal Threshold:", optimal_threshold)
23 print("Maximum F1 Score:", optimal_f1)
24
25 # Convert reconstruction error to binary predictions
26 binary_predictions = [1 if error > optimal_threshold else 0 for error in r
27
28 # Calculate F1 score
29 final_f1_score = f1_score(y_test_final, binary_predictions)
30
31 print("F1 score:", final_f1_score)
32

```

7/7 ————— 0s 833us/step

Optimal Threshold: 1.0259181923459462

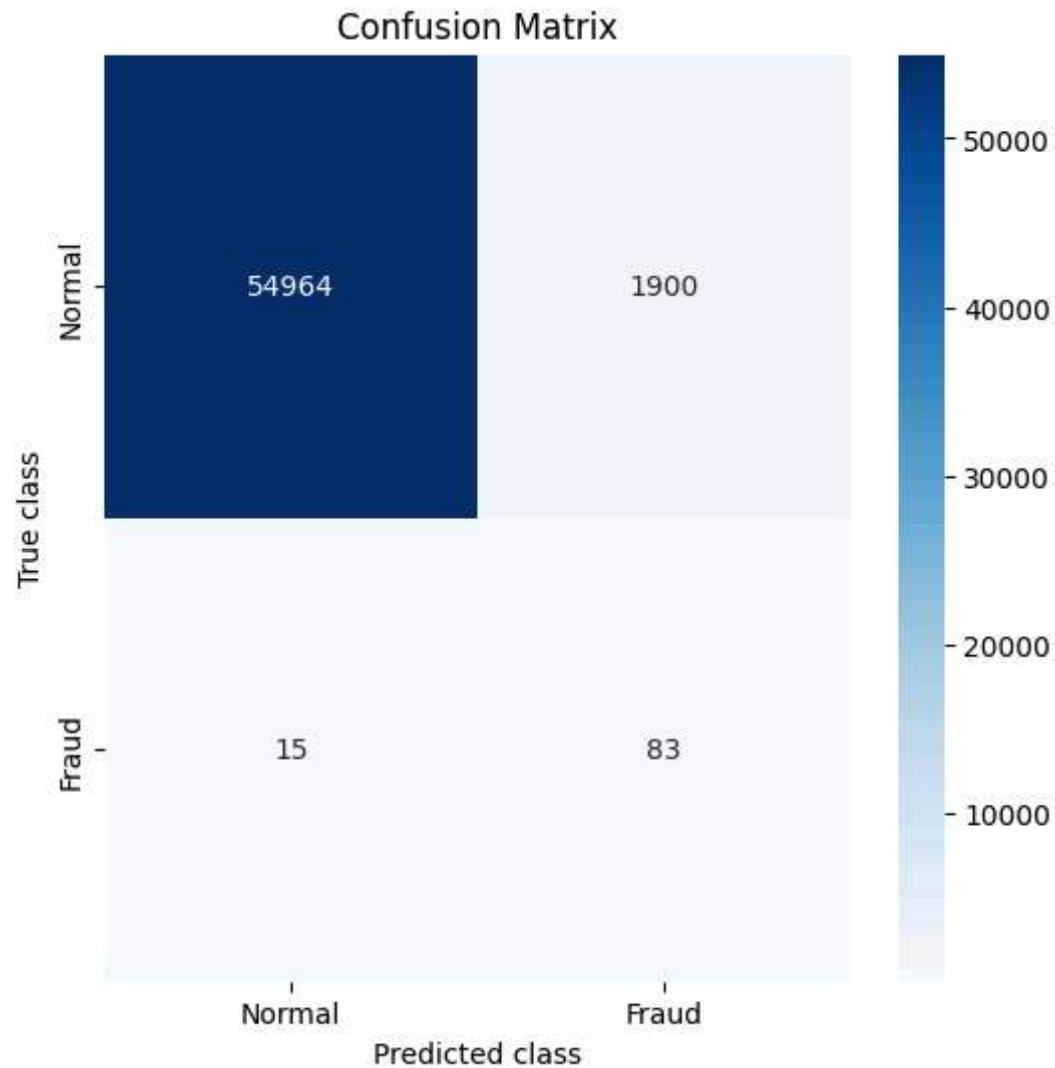
Maximum F1 Score: 0.9175257731958762

F1 score: 0.9119170984455959

Confusion matrix

This confusion matrix shows that for a binary classification problem, the classifier correctly identified 54,964 normal cases and 83 fraud cases. However, it incorrectly classified 1,900 normal cases as fraud (false positives) and failed to identify 15 fraud cases (false negatives). The matrix suggests the model is relatively good at detecting normal cases but less effective at identifying fraud.

```
In [80]:  
1 LABELS = ["Normal", "Fraud"]  
2  
3 # Confusion Matrix  
4 prediction = [1 if e > optimal_threshold else 0 for e in mse]  
5 conf_matrix = confusion_matrix(y_test, prediction)  
6  
7 plt.figure(figsize=(6, 6))  
8 sns.heatmap(conf_matrix, xticklabels=LABELS, yticklabels=LABELS, annot=True)  
9 plt.title("Confusion Matrix")  
10 plt.ylabel('True class')  
11 plt.xlabel('Predicted class')  
12 plt.show()
```



VAE

Architecture

Training

```
In [ ]:
 1 class Sampling(layers.Layer):
 2     """Uses (z_mean, z_log_var) to sample z, the vector encoding a digit."""
 3     def call(self, inputs):
 4         z_mean, z_log_var = inputs
 5         batch = tf.shape(z_mean)[0]
 6         dim = tf.shape(z_mean)[1]
 7         epsilon = K.random_normal(shape=(batch, dim))
 8         return z_mean + tf.exp(0.5 * z_log_var) * epsilon
 9
10 # Define encoder model.
11 input_dim = X_train.shape[1]
12 latent_dim = 2
13
14 encoder_inputs = layers.Input(shape=(input_dim,))
15 x = layers.Dense(64, activation='relu')(encoder_inputs) # Increased to 64
16 x = layers.Dense(128, activation='relu')(x) # Added another layer
17 z_mean = layers.Dense(latent_dim, name='z_mean')(x)
18 z_log_var = layers.Dense(latent_dim, name='z_log_var')(x)
19 z = Sampling()([z_mean, z_log_var])
20 encoder = models.Model(encoder_inputs, [z_mean, z_log_var, z], name='encoder')
21
22 # Increase the width and depth of the decoder as well.
23 latent_inputs = layers.Input(shape=(latent_dim,), name='z_sampling')
24 x = layers.Dense(128, activation='relu')(latent_inputs) # Matched with encoder
25 x = layers.Dense(64, activation='relu')(x) # Matched with encoder
26 decoder_outputs = layers.Dense(input_dim, activation='sigmoid')(x)
27 decoder = models.Model(latent_inputs, decoder_outputs, name='decoder')
28
29 # Define VAE model.
30 class VAE(models.Model):
31     def __init__(self, encoder, decoder, **kwargs):
32         super(VAE, self).__init__(**kwargs)
33         self.encoder = encoder
34         self.decoder = decoder
35
36     def call(self, inputs):
37         z_mean, z_log_var, z = self.encoder(inputs)
38         reconstructed = self.decoder(z)
39         # Add KL divergence regularization Loss.
40         kl_loss = -0.5 * K.sum(1 + z_log_var - K.square(z_mean) - K.exp(z_mean))
41         kl_loss = K.mean(kl_loss) # Ensure the KL loss is scalar.
42         self.add_loss(kl_loss)
43         return reconstructed
44 vae = VAE(encoder, decoder)
45 optimizer = Adam(learning_rate=1e-3)
46
47
48 # Compile VAE
49 vae = VAE(encoder, decoder)
50 optimizer = Adam(learning_rate=1e-3)
51
52 # Compile VAE
53 vae.compile(optimizer=optimizer, loss=MeanSquaredError())
54
55 # Prepare training and validation data.
56 X_train_np = X_train.to_numpy().astype(np.float32)
57 X_test_np = X_test.to_numpy().astype(np.float32)
```

```
58
59 # Ensure data is normalized if it's not already
60 X_train_np = X_train_np / np.max(X_train_np)
61 X_test_np = X_test_np / np.max(X_test_np)
62
63 # Training Loop.
64 batch_size = 1000
65 epochs = 50 # Increase the number of epochs
66
67 # Add callbacks for early stopping and Learning rate reduction
68 callbacks = [
69     EarlyStopping(monitor='val_loss', patience=10, verbose=1, restore_best=True),
70     ReduceLROnPlateau(monitor='val_loss', factor=0.5, patience=5, verbose=1)
71 ]
72
73 # Fit the model using data.
74 history = vae.fit(
75     X_train_np, X_train_np,
76     epochs=epochs,
77     batch_size=batch_size,
78     validation_data=(X_test_np, X_test_np),
79     callbacks=callbacks
80 )
81
```

```
Epoch 1/50
228/228 2s 3ms/step - loss: 0.0969 - val_loss: 7.3677e-04 - learning_rate: 0.0010
Epoch 2/50
228/228 1s 3ms/step - loss: 1.8328e-04 - val_loss: 6.0824e-04 - learning_rate: 0.0010
Epoch 3/50
228/228 1s 2ms/step - loss: 1.0208e-04 - val_loss: 5.8690e-04 - learning_rate: 0.0010
Epoch 4/50
228/228 1s 3ms/step - loss: 8.7061e-05 - val_loss: 5.7919e-04 - learning_rate: 0.0010
Epoch 5/50
228/228 1s 2ms/step - loss: 8.2661e-05 - val_loss: 5.7548e-04 - learning_rate: 0.0010
Epoch 6/50
228/228 1s 2ms/step - loss: 8.0003e-05 - val_loss: 5.7337e-04 - learning_rate: 0.0010
Epoch 7/50
228/228 1s 2ms/step - loss: 7.8655e-05 - val_loss: 5.7126e-04
```

In []: 1 encoder.summary()

Model: "encoder"

Layer (type)	Output Shape	Param #	Connected to
input_layer_2 (InputLayer)	(None, 29)	0	-
dense_8 (Dense)	(None, 64)	1,920	input_layer_2[0]...
dense_9 (Dense)	(None, 128)	8,320	dense_8[0][0]
z_mean (Dense)	(None, 2)	258	dense_9[0][0]
z_log_var (Dense)	(None, 2)	258	dense_9[0][0]
sampling (Sampling)	(None, 2)	0	z_mean[0][0], z_log_var[0][0]

Total params: 10,756 (42.02 KB)

Trainable params: 10,756 (42.02 KB)

Non-trainable params: 0 (0.00 B)

In []: 1 decoder.summary()

Model: "decoder"

Layer (type)	Output Shape	Param #
z_sampling (InputLayer)	(None, 2)	0
dense_10 (Dense)	(None, 128)	384
dense_11 (Dense)	(None, 64)	8,256
dense_12 (Dense)	(None, 29)	1,885

Total params: 10,525 (41.11 KB)

Trainable params: 10,525 (41.11 KB)

Non-trainable params: 0 (0.00 B)

```
In [ ]: 1 vae.summary()
```

Model: "vae_1"

Layer (type)	Output Shape	Param #
encoder (Functional)	?	10,756
decoder (Functional)	?	10,525

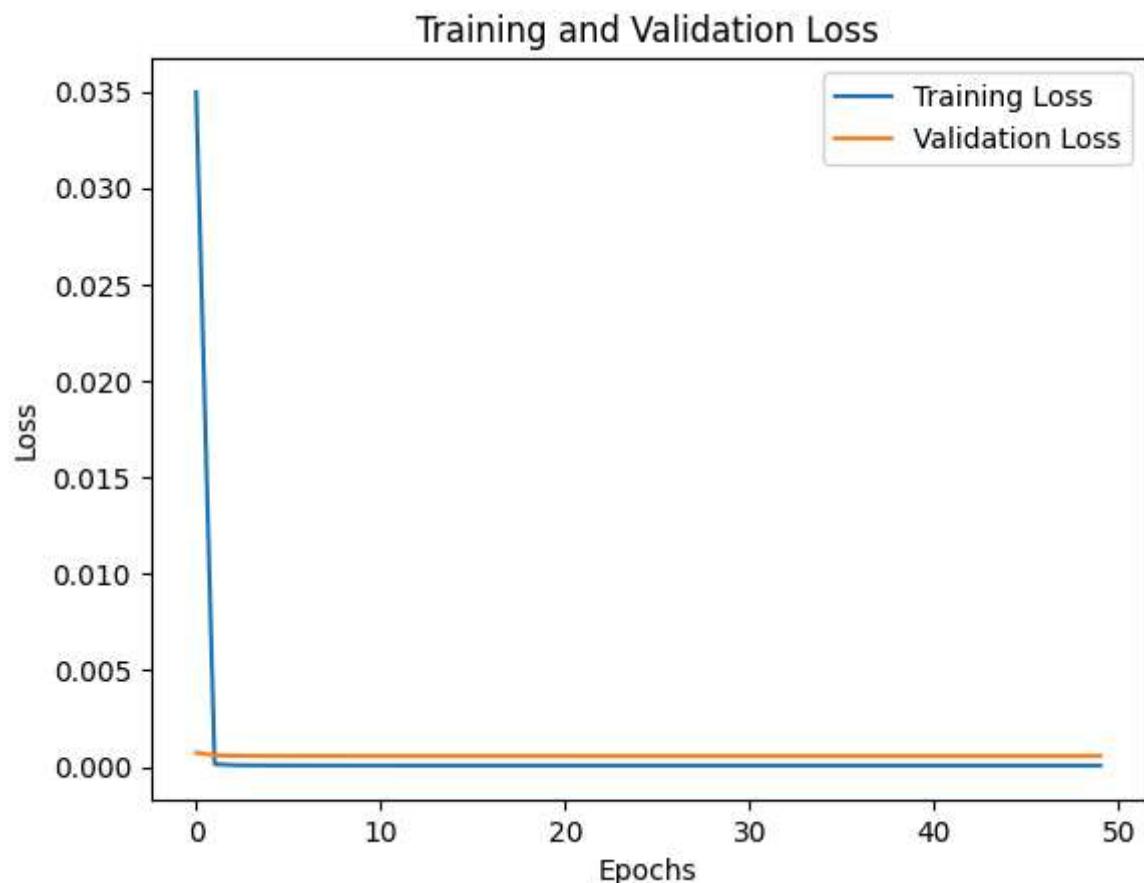
Total params: 63,845 (249.40 KB)

Trainable params: 21,281 (83.13 KB)

Non-trainable params: 0 (0.00 B)

Optimizer params: 42,564 (166.27 KB)

```
In [ ]: 1 # Plotting training Loss  
2 plt.plot(history.history['loss'], label='Training Loss')  
3 # Plotting validation Loss  
4 plt.plot(history.history['val_loss'], label='Validation Loss')  
5  
6 # Adding title and Labels  
7 plt.title('Training and Validation Loss')  
8 plt.xlabel('Epochs')  
9 plt.ylabel('Loss')  
10 plt.legend()  
11  
12 # Display the plot  
13 plt.show()  
14
```



```
In [ ]: 1 def reconstruction_log_prob(eval_samples, reconstruct_samples_n):  
2     # Get z_mean and z_log_var from the encoder  
3     z_mean, z_log_var, _ = encoder(eval_samples)  
4  
5     # Prepare to sample from the latent space  
6     epsilon = tf.random.normal(shape=(reconstruct_samples_n, z_mean.shape[1]))  
7     z_samples = z_mean[tf.newaxis, :, :] + tf.exp(0.5 * z_log_var[tf.newaxis, :, :]) * epsilon  
8  
9     # Decode each sample  
10    decoded_samples = tf.map_fn(decoder, z_samples, dtype=tf.float32) # A  
11  
12    # Assuming the output of decoder is the reconstruction itself, we use  
13    # Here we manually compute MSE as a proxy for log_prob under a Gaussian  
14    mse = tf.reduce_mean(tf.square(decoded_samples - eval_samples[tf.newaxis, :, :]))  
15    log_prob = -mse # Negate MSE to simulate log probability (higher is better)  
16  
17    return tf.reduce_mean(log_prob, axis=0) # Average over all samples  
18
```

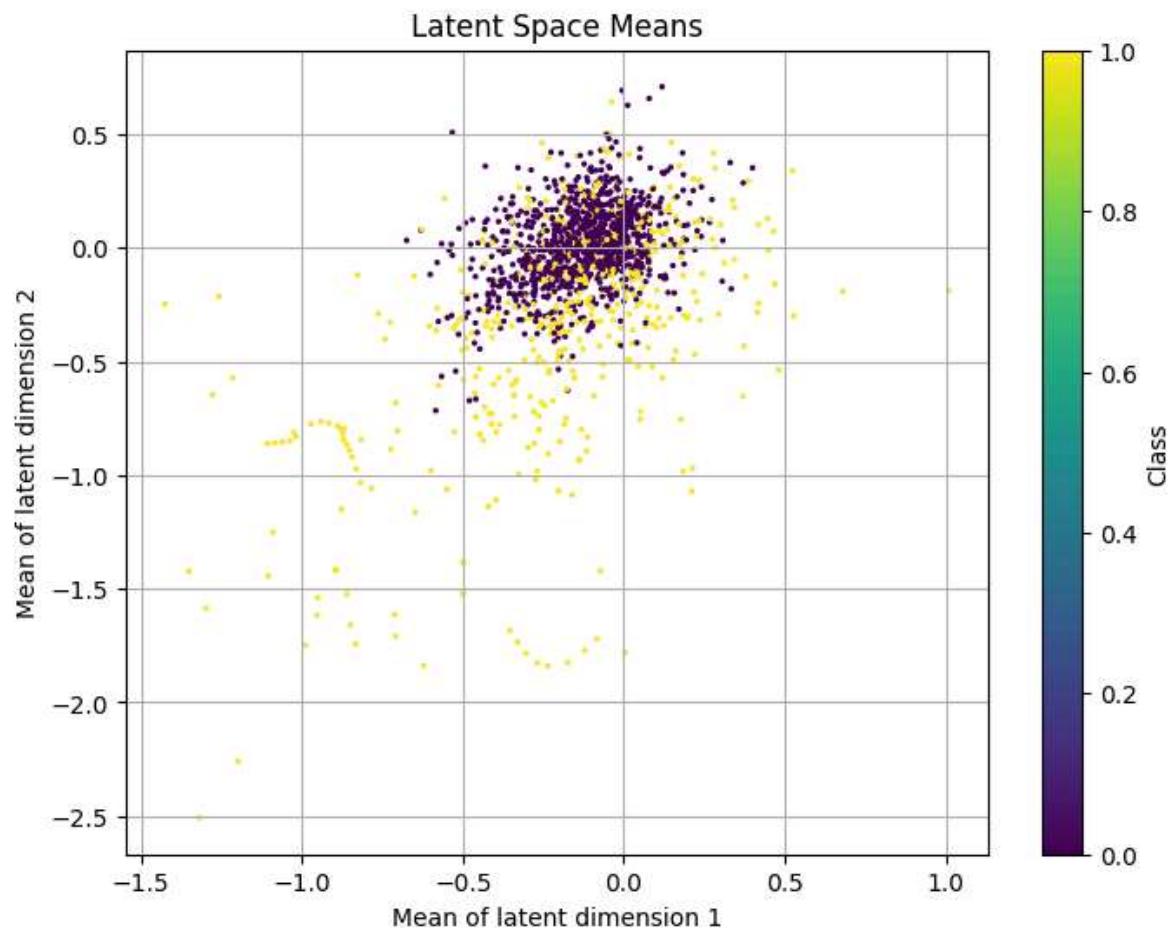
Results

Latent space mean

This scatter plot represents data points in a two-dimensional latent space, with each point colored according to its class. The dense cluster of purple points near the center suggests a concentration of one class, while the sparser yellow points indicate another class. The plot could be showing how a model like a Variational Autoencoder separates features of different classes in its learned latent space

```
In [ ]: 1 # Use the encoder to get the mean, Log variance, and sampled z outputs
2 z_mean, z_log_var, _ = encoder.predict(X) # Use the predict method and call it z_mean
3
4 # Plot these means colored by their corresponding label `Y`
5 plt.figure(figsize=(8, 6))
6 scatter = plt.scatter(z_mean[:, 0], z_mean[:, 1], c=Y, cmap='viridis', s=2)
7 plt.colorbar(scatter, label='Class')
8 plt.title('Latent Space Means')
9 plt.xlabel('Mean of latent dimension 1')
10 plt.ylabel('Mean of latent dimension 2')
11 plt.grid(True)
12 plt.show()
13
```

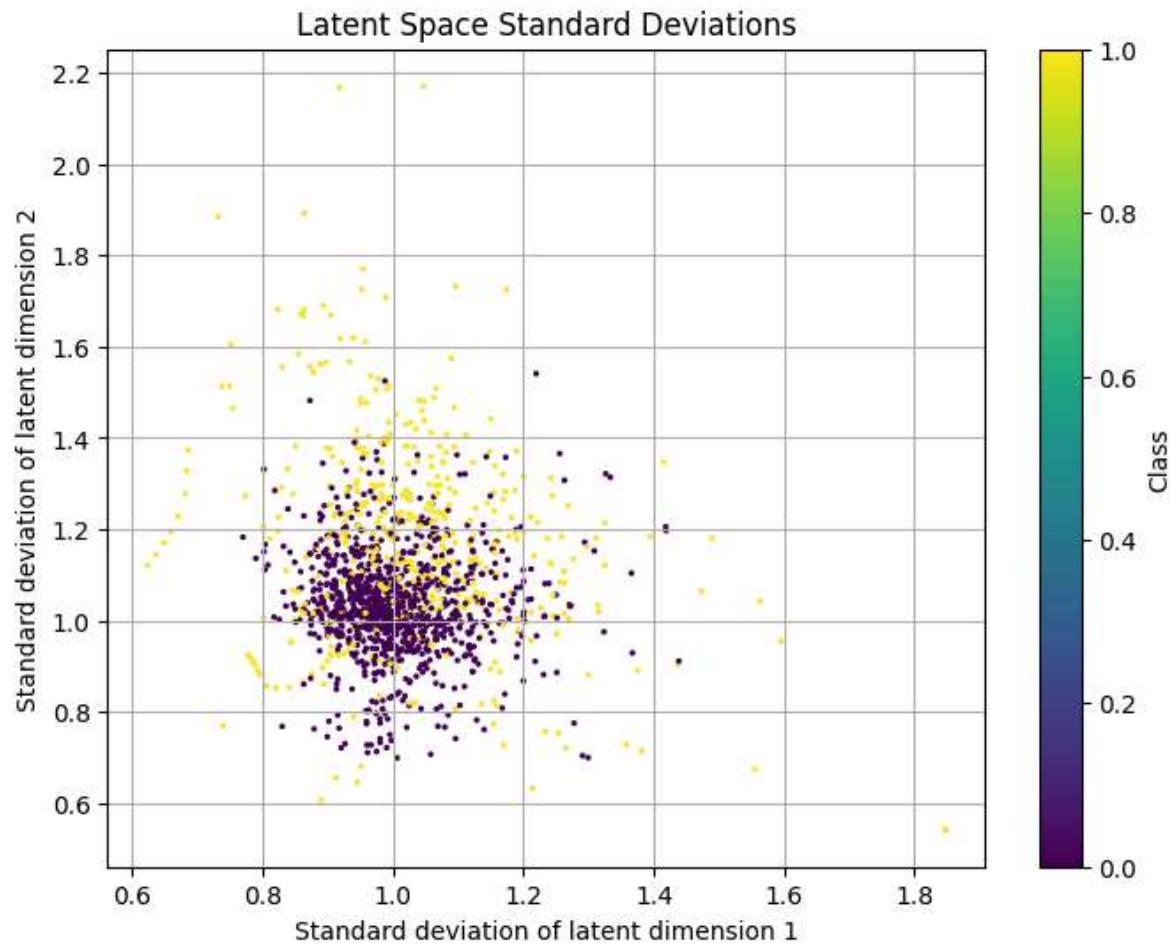
47/47 0s 1ms/step



Latent space standard deviation

The scatter plot illustrates data points in terms of their standard deviation across two latent dimensions, colored by class. A cluster of purple points indicates a class with lower variability, while yellow points dispersed throughout indicate higher variability in another class. This visualization can be useful for understanding the spread and overlap of the different classes in a model's latent space

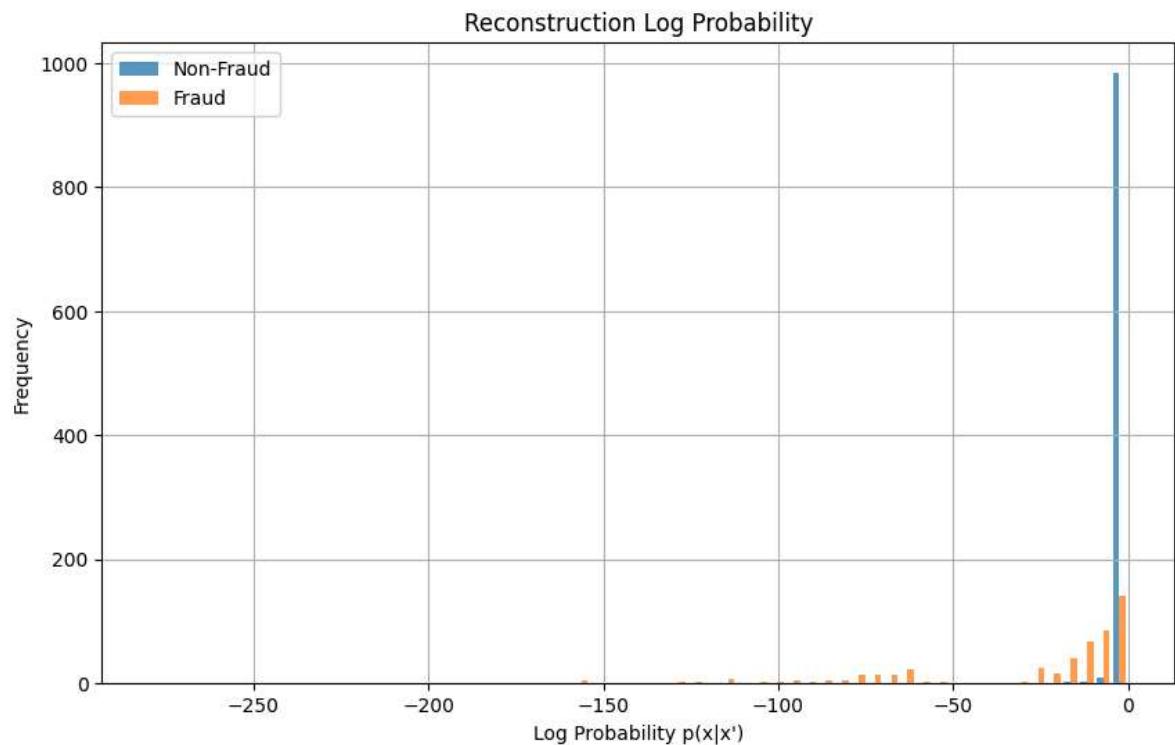
```
In [ ]: 1 latent_x_std = tf.exp(0.5 * z_log_var)
2
3 plt.figure(figsize=(8, 6))
4 scatter = plt.scatter(latent_x_std[:, 0], latent_x_std[:, 1], c=Y, cmap='viridis')
5 plt.colorbar(scatter, label='Class')
6 plt.title('Latent Space Standard Deviations')
7 plt.xlabel('Standard deviation of latent dimension 1')
8 plt.ylabel('Standard deviation of latent dimension 2')
9 plt.grid(True)
10 plt.show()
```



Reconstruction log probability

The histogram compares the log probability of reconstruction for non-fraud and fraud transactions. Non-fraud transactions show a high frequency of high log probability scores, indicating the model reconstructs them with high likelihood. Fraud transactions, however, have a lower frequency and are spread out across lower log probability values, suggesting they are less likely to be reconstructed accurately by the model

```
In [ ]: 1 # Compute Log probabilities
2 reconstruct_samples_n = 50 # Example: 100 samples
3 x_log_prob = reconstruction_log_prob(X, reconstruct_samples_n)
4
5 # Create a histogram of Log probabilities for both classes
6 plt.figure(figsize=(10, 6))
7 plt.hist([x_log_prob[Y == 0], x_log_prob[Y == 1]], bins=60, label=['Non-Fraud', 'Fraud'])
8 plt.title('Reconstruction Log Probability')
9 plt.xlabel("Log Probability p(x|x')")
```



```
In [ ]: 1 LABELS = ["Normal", "Fraud"]
2 # Predict the reconstruction loss on the test set
3 test_predictions = vae.predict(X_test)
4 mse = np.mean(np.power(X_test - test_predictions, 2), axis=1)
5
6 # Create a DataFrame with true class Labels and the calculated MSE
7 error_df = pd.DataFrame({'reconstruction_error': mse,
8                           'true_class': y_test})
```

1781/1781 ━━━━━━━━ 2s 1ms/step

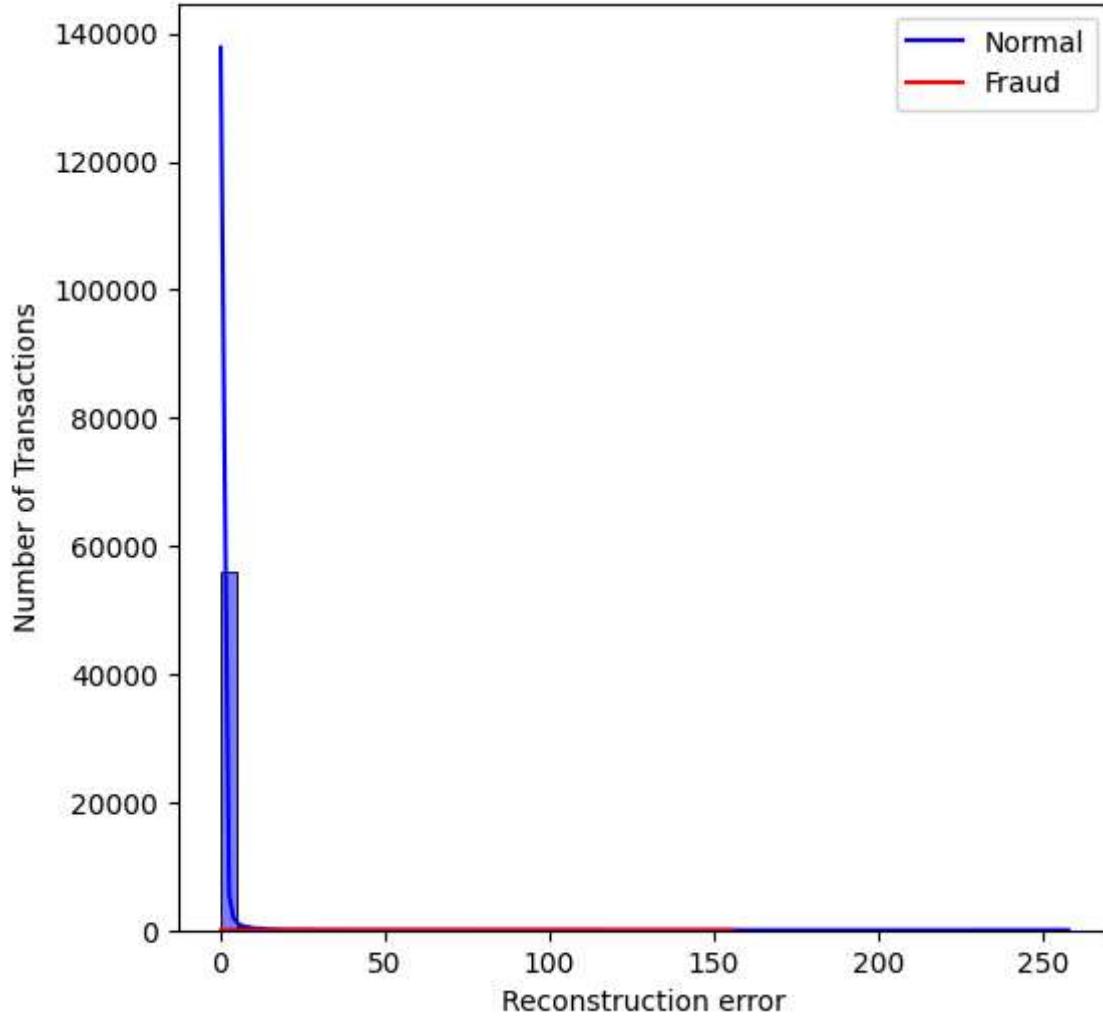
Reconstruction Error for Fraud Transactions

Normal transactions are tightly clustered around a low reconstruction error, indicating accurate model predictions. Fraud transactions are less frequent and have a wider distribution of reconstruction errors, likely due to their anomalous nature which the model finds harder to

reconstruct accurately

In []:

```
1 # Visualizing the Reconstruction Error for different classes
2 fig, ax = plt.subplots(figsize=(6, 6))
3 normal_error_df = error_df[(error_df['true_class'] == 0)]
4 fraud_error_df = error_df[error_df['true_class'] == 1]
5
6 sns.histplot(normal_error_df['reconstruction_error'], bins=50, kde=True, color='blue')
7 sns.histplot(fraud_error_df['reconstruction_error'], bins=50, kde=True, color='red')
8 plt.xlabel('Reconstruction error')
9 plt.ylabel('Number of Transactions')
10 plt.legend(['Normal', 'Fraud'])
11 plt.show()
```

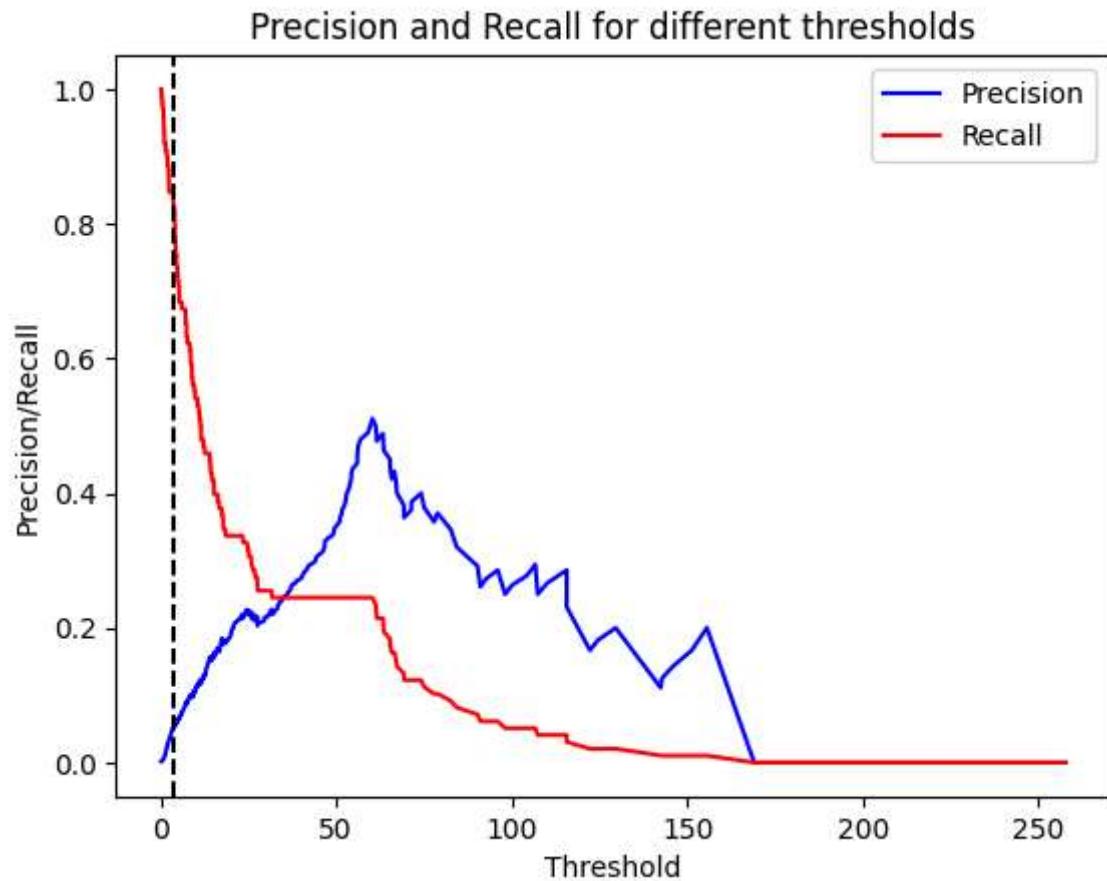


Recall vs Precision

The graph shows precision and recall values for various thresholds in a classification model. Precision starts high but quickly drops and fluctuates as the threshold increases, while recall decreases more steadily. The dashed line may represent an optimal threshold where a balance between precision and recall is achieved, often used to optimize a model's performance on both metrics.

```
In [ ]: 1 # Determine the optimal threshold
2 precision, recall, threshold = precision_recall_curve(error_df.true_class,
3 #threshold_f1 = threshold[np.argmax(2 * recall[:-1] * precision[:-1]) / (recall[:-1] + precision[:-1])])
4
5 threshold_f1 = 3.8
6
7 print(threshold_f1)
8
9 plt.plot(threshold, precision[:-1], 'b', label='Precision')
10 plt.plot(threshold, recall[:-1], 'r', label='Recall')
11 plt.axvline(x=threshold_f1, color='k', linestyle='--')
12 plt.title('Precision and Recall for different thresholds')
13 plt.xlabel('Threshold')
14 plt.ylabel('Precision/Recall')
15 plt.legend()
16 plt.show()
```

3.8

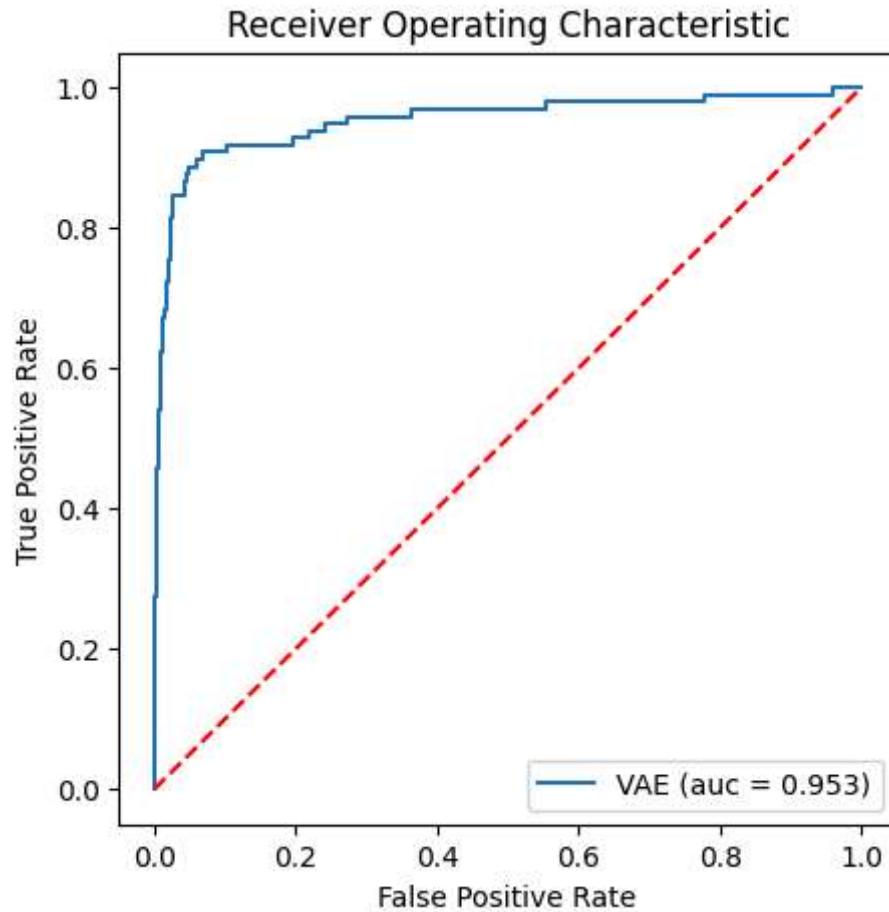


ROC curve

The ROC curve depicts the performance of a Variational Autoencoder (VAE) with an AUC of 0.953, which is quite high, indicating strong discriminative power. The curve stays close to the top-left corner, suggesting a high true positive rate with a low false positive rate, which is desirable in many classification tasks.

In []:

```
1 # ROC Curve
2 fpr, tpr, thresholds = roc_curve(error_df.true_class, error_df.reconstruct)
3 roc_auc = auc(fpr, tpr)
4
5 plt.figure(figsize=(5, 5), dpi=100)
6 plt.plot(fpr, tpr, linestyle='-', label='VAE (auc = %0.3f)' % roc_auc)
7 plt.plot([0, 1], [0, 1], 'r--')
8 plt.xlabel('False Positive Rate')
9 plt.ylabel('True Positive Rate')
10 plt.title('Receiver Operating Characteristic')
11 plt.legend()
12 plt.show()
```

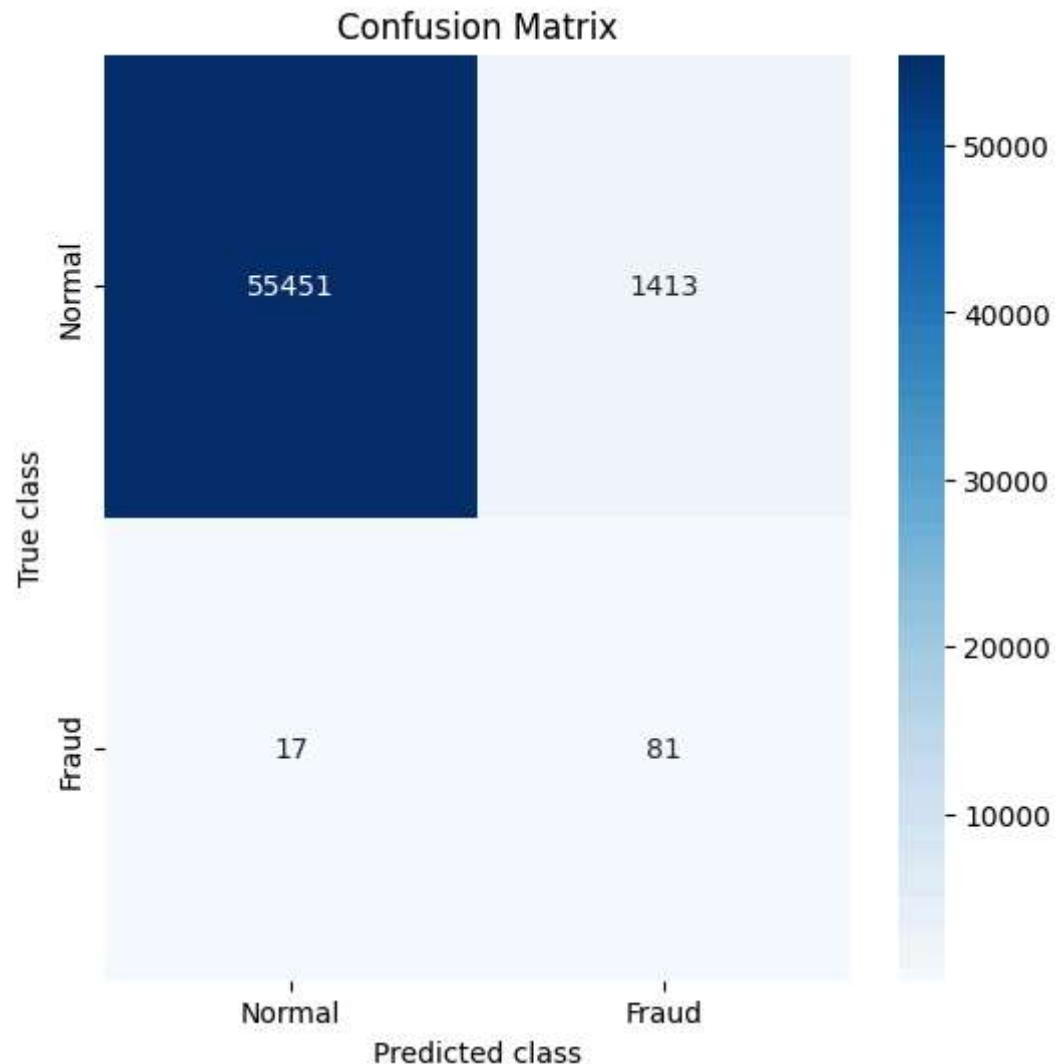


Confusion matrix

The confusion matrix shows a classifier's predictions, with 55,451 true negatives and 81 true positives, indicating correct classifications for normal and fraud cases respectively. However, there are 1,413 false positives and 17 false negatives, showing instances where the model incorrectly classified normal transactions as fraud and missed identifying some fraud transactions.

In []:

```
1 # Confusion Matrix
2 prediction = [1 if e > threshold_f1 else 0 for e in error_df.reconstruction_error]
3 conf_matrix = confusion_matrix(error_df.true_class, prediction)
4
5 plt.figure(figsize=(6, 6))
6 sns.heatmap(conf_matrix, xticklabels=LABELS, yticklabels=LABELS, annot=True)
7 plt.title("Confusion Matrix")
8 plt.ylabel('True class')
9 plt.xlabel('Predicted class')
10 plt.show()
```



```
In [ ]: 1 reconstructed = vae.predict(X_test_np)
2 reconstruction_error = np.mean(np.square(X_test_np - reconstructed), axis=-1)
3
4 # Calculate optimal threshold and F1 score using reconstruction error
5 def calculate_optimal_f1(y_true, errors):
6     precision, recall, thresholds = precision_recall_curve(y_true, errors)
7     f1_scores = 2 * (precision * recall) / (precision + recall)
8     f1_scores = np.nan_to_num(f1_scores) # Handle potential NaNs
9     max_index = np.argmax(f1_scores)
10    return f1_scores[max_index], thresholds[max_index]
11
```

1781/1781 ————— 1s 548us/step

```
In [ ]: 1 # Function to balance the dataset
2 def balance_classes(X, y):
3     # Separate the classes
4     class_0 = X[y == 0]
5     class_1 = X[y == 1]
6
7     # Find the smaller class size
8     n_samples = min(len(class_0), len(class_1))
9
10    # Resample classes to the size of the smaller class
11    class_0_downsampled = resample(class_0, replace=False, n_samples=n_samples)
12    class_1_downsampled = resample(class_1, replace=False, n_samples=n_samples)
13
14    # Concatenate back to a single array
15    X_balanced = np.vstack((class_0_downsampled, class_1_downsampled))
16    y_balanced = np.array([0] * n_samples + [1] * n_samples)
17
18    return X_balanced, y_balanced
19
20 # Balancing the classes in the test set
21 X_test_balanced, y_test_balanced = balance_classes(X_test_np, y_test)
22
23 # Prediction using the VAE
24 reconstructed_balanced = vae.predict(X_test_balanced)
25 reconstruction_error_balanced = np.mean(np.square(X_test_balanced - reconstructed_balanced))
26
27 # Calculate optimal threshold and F1 score using balanced reconstruction error
28 optimal_f1, optimal_threshold = calculate_optimal_f1(y_test_balanced, reconstruction_error_balanced)
29 print("Optimal Threshold:", optimal_threshold)
30 print("Optimal F1 Score:", optimal_f1)
```

7/7 ————— 0s 656us/step

Optimal Threshold: 0.0007648746

Optimal F1 Score: 0.923076923076923