

# Fraud Analytics (CS6890)

**Assignment:**

**5**

**Title :** Synthetic data generation using Variational Autoencoder

**Team Details :**

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

## Exploratory Data Analysis

### Importing Required Modules

In [103]:

```
1 import pandas as pd
2 from sklearn.preprocessing import LabelEncoder
3 import numpy as np
4 import pandas as pd
5 from tensorflow.keras.models import Model, load_model
6 from tensorflow.keras.layers import Input, Dense, Lambda, Concatenate, Dropout
7 from tensorflow.keras.optimizers import Adam
8 from tensorflow.keras.callbacks import EarlyStopping, LearningRateScheduler
9 from tensorflow.keras import backend as K
10 import pandas as pd
11 from sklearn.model_selection import train_test_split
12 from sklearn.preprocessing import StandardScaler, MinMaxScaler
13 from tensorflow.keras.losses import categorical_crossentropy
14 import matplotlib.pyplot as plt
```

## Checking for duplicate rows

```
In [102]: 1 credit_card_data = pd.read_csv('card_transaction.v1.csv')
2 credit_card_data[credit_card_data.duplicated()]
```

Out[102]:

	User	Card	Year	Month	Day	Time	Amount	Use Chip	Merchant Name	Merchant City	Merchant State
193878	13	0	2015	4	28	396	100.0	0	4487	725	71 4
195221	13	1	2006	7	19	454	100.0	2	4487	725	71 4
195330	13	1	2006	10	15	431	100.0	2	4487	725	71 4
244172	17	2	2011	12	12	986	120.0	2	4487	5722	3 8
1360769	109	2	1999	12	9	256	120.0	2	4487	2919	50 4
1373587	109	2	2019	3	8	262	120.0	0	4487	2919	50 4
2469536	205	2	2012	2	8	1053	82.0	2	12586	3305	91 1

## Checking for Missing/ Null Values

```
In [64]: 1 credit_card_data.isnull().sum(axis=0)
```

```
Out[64]: User          0
Card          0
Year          0
Month         0
Day           0
Time          0
Amount        1
Use Chip      1
Merchant Name 1
Merchant City 1
Merchant State 321041
Zip           338205
MCC           1
Errors?       2515428
Is Fraud?     1
dtype: int64
```

```
In [65]: 1 credit_card_data[credit_card_data['Merchant State'].isnull()]['Merchant C:
```

```
Out[65]: Merchant City
          ONLINE    321040
          Name: count, dtype: int64
```

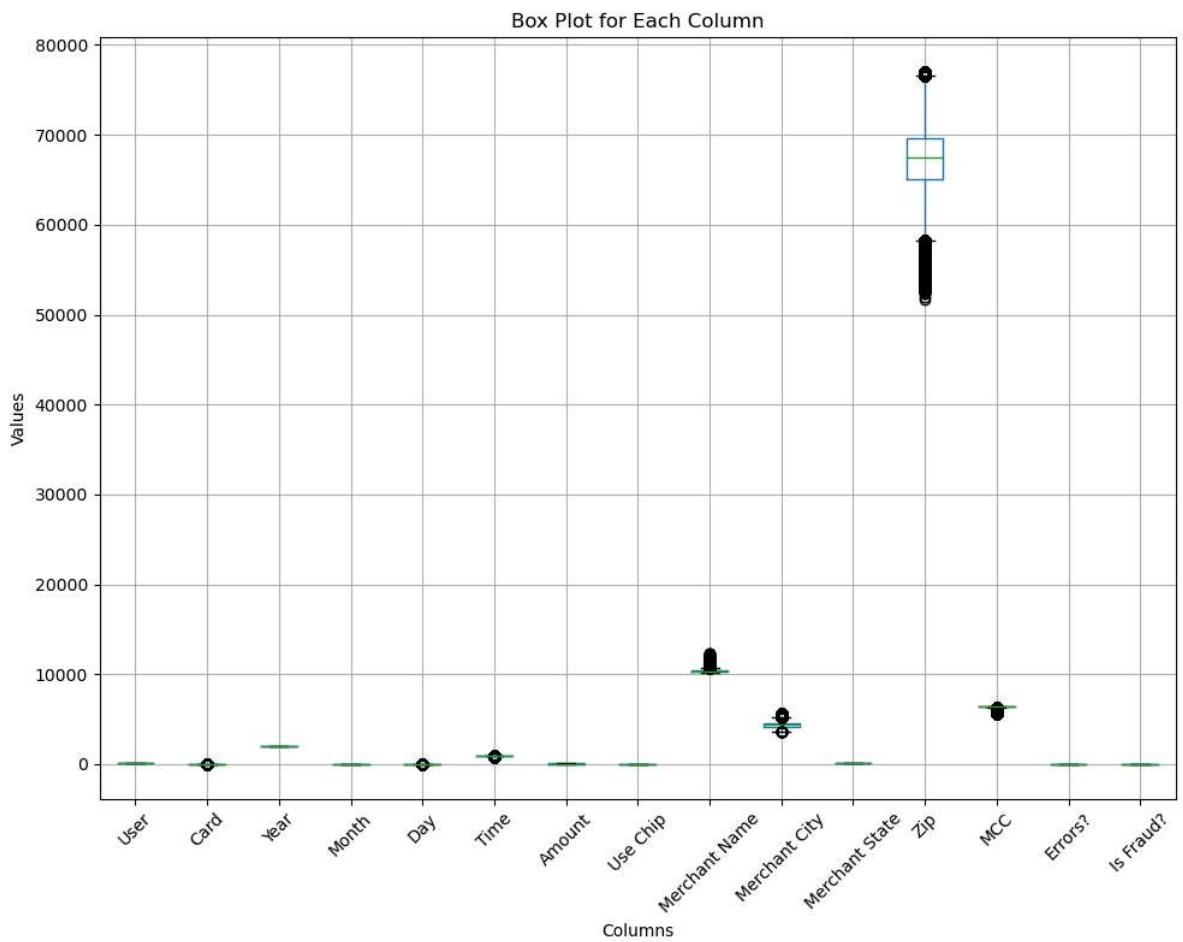
```
In [66]: 1 credit_card_data[credit_card_data['Merchant City'] == ' ONLINE']['Merchant C:
```

```
Out[66]: Series([], Name: count, dtype: int64)
```

## Plots

In [108]:

```
1 import pandas as pd
2 import matplotlib.pyplot as plt
3
4 # Assuming df is your DataFrame containing the dataset
5 # Replace df with your actual DataFrame name
6
7 # Define the columns for which you want to create box plots
8 columns_to_plot = ['User', 'Card', 'Year', 'Month', 'Day', 'Time', 'Amount',
9                     'Merchant Name', 'Merchant City', 'Merchant State', 'Zip',
10                    'Use Chip', 'MCC', 'Errors?', 'Is Fraud?']
11
12 # Create a box plot for each column
13 plt.figure(figsize=(10, 8))
14 df[columns_to_plot].boxplot()
15 plt.xticks(rotation=45)
16 plt.title('Box Plot for Each Column')
17 plt.xlabel('Columns')
18 plt.ylabel('Values')
19 plt.tight_layout()
20 plt.show()
```



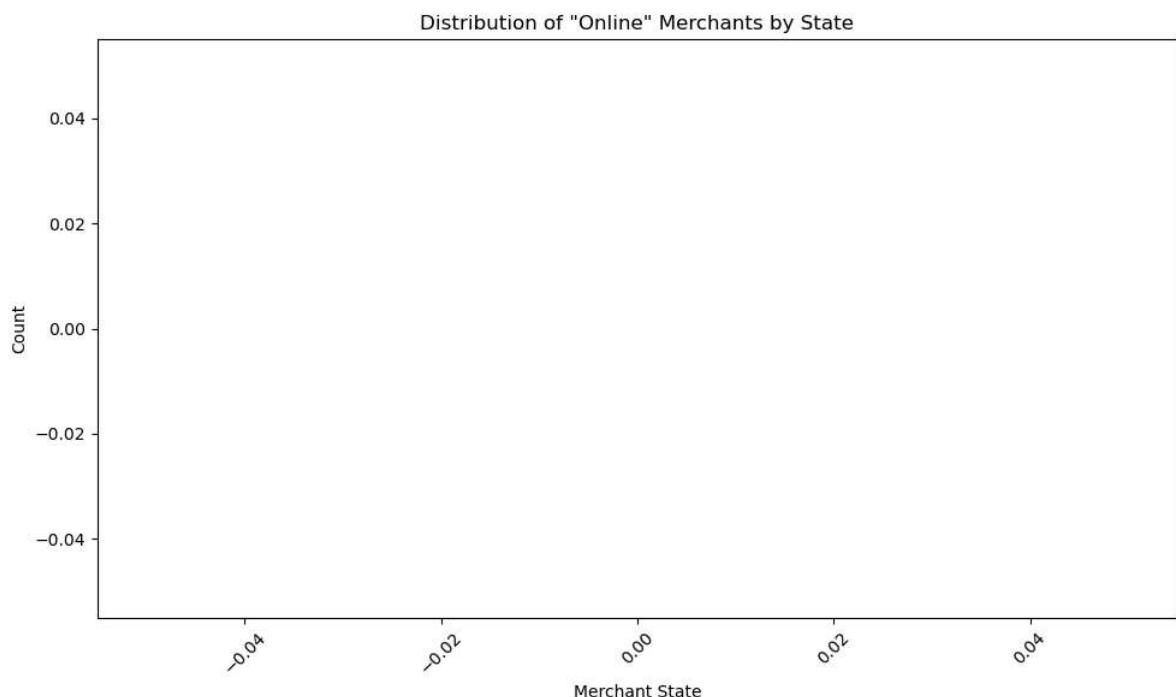
## Finding relation between the missing values and the observed data

### 'Merchant City' Vs 'Merchnat State'

- Missing values can be classified into three types:
  - Missing at Random (MAR),
  - Missing Completely at Random (MCAR), and
  - Missing Not at Random (MNAR)
- Upon analyzing our dataset, we identified that the missing values in the "Merchant City" column fall under the category of **Missing at Random (MAR)**.
- Specifically, for the ' ONLINE' category in Merchant City, all corresponding values in the Merchant State column were missing. To address this, we utilized a random constant imputation technique. This involved assigning a random constant value to all missing entries within the 'ONLINE' category of Merchant City. By doing so, we were able to uniformly impute the missing values in a consistent manner across this category.

In [107]:

```
1 filtered_df = credit_card_data[credit_card_data['Merchant City'] == ' ONLINE']
2
3 # Plotting the graph
4 plt.figure(figsize=(10, 6))
5 plt.bar(filtered_df['Merchant State'], filtered_df['Merchant City'].count())
6 plt.xlabel('Merchant State')
7 plt.ylabel('Count')
8 plt.title('Distribution of "Online" Merchants by State')
9 plt.xticks(rotation=45)
10 plt.tight_layout()
11 plt.show()
```



## 'Merchant City' vs 'Zip'

- In our data analysis, we encountered a similar scenario when examining the "Merchant City" and "Zip" columns. It was evident that the absence of values in the "Zip" column correlated with various categories within the "Merchant City" column i.e **Missing as a Random (MAR) type**. Specifically, certain cities, including the 'ONLINE' category, exhibited missing values in their respective "Zip" entries.
- To tackle this issue, we first identified all unique categories within "Merchant City" that had missing "Zip" values. Subsequently, we assigned a distinct constant value to each of these identified categories. This approach aimed to establish a clear correlation between missing "Zip" values and specific "Merchant City" categories. We then utilized these unique constants to fill in the missing "Zip" entries, ensuring a consistent and systematic imputation method across different categories.

```
In [67]: 1 credit_card_data[credit_card_data['Zip'].isnull()]['Merchant City'].value_
```

```
Out[67]: Merchant City
    ONLINE      321040
    Cancun       1973
    Rome         1247
    Mexico City   945
    London        791
    ...
    Vaduz          3
    Palikir        3
    Bratislava     2
    Apia           1
    Freetown        1
Name: count, Length: 114, dtype: int64
```

```
In [68]: 1 credit_card_data[credit_card_data['Merchant City'] == ' ONLINE']['Zip'].va
```

```
Out[68]: Series([], Name: count, dtype: int64)
```

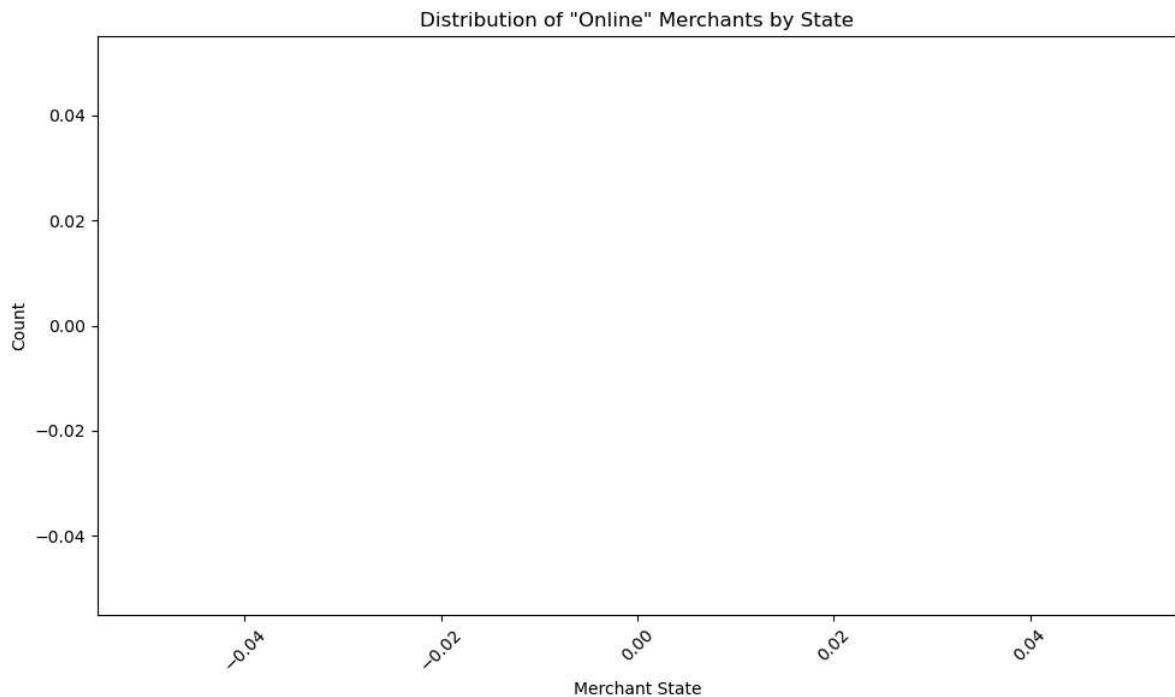
```
In [69]: 1 missing_zip_categories = credit_card_data[credit_card_data['Zip'].isna()]|
```

```
In [70]: 1 missing_zip_categories.tolist()
```

```
Out[70]: ['ONLINE',
'Zurich',
'Tallinn',
'Tokyo',
'Cabo San Lucas',
'Cancun',
'Santo Domingo',
'Beijing',
'Lisbon',
'Rome',
'Manila',
'Kingston',
'Shanghai',
'Toronto',
'Paris',
'Oslo',
'Geneva',
'Wellington',
'Amsterdam',
'-' ]
```

```
In [71]: 1 element_dict = {}
2 for index, element in enumerate(missing_zip_categories, start=1):
3     element_dict[element] = index
```

```
In [109]: 1 filtered_df1 = credit_card_data[credit_card_data['Zip'].isnull()]
2
3 # Plotting the graph
4 plt.figure(figsize=(10, 6))
5 plt.bar(filtered_df1['Zip'], filtered_df1['Merchant City'].count())
6 plt.xlabel('Merchant State')
7 plt.ylabel('Count')
8 plt.title('Distribution of "Online" Merchants by State')
9 plt.xticks(rotation=45)
10 plt.tight_layout()
11 plt.show()
```



## Imputing Missing Values

- The missing values in the 'Amount', 'Use Chip', 'Merchant Name', 'Merchant City', 'MCC', and 'Is Fraud?' columns were addressed by employing a mode imputation technique. Since each of these columns had only one missing value and they were categorical in nature, filling the missing entries with the mode (most frequent value) of each column ensured that the dataset remained consistent and preserved the categorical information effectively.

```
In [72]: 1 credit_card_data['Amount'].fillna(credit_card_data['Amount'].mode()[0], inplace=True)
2 credit_card_data['Use Chip'].fillna(credit_card_data['Use Chip'].mode()[0], inplace=True)
3 credit_card_data['Merchant Name'].fillna(credit_card_data['Merchant Name'].mode()[0], inplace=True)
4 credit_card_data['Merchant City'].fillna(credit_card_data['Merchant City'].mode()[0], inplace=True)
5 credit_card_data['MCC'].fillna(credit_card_data['MCC'].mode()[0], inplace=True)
6 credit_card_data['Is Fraud?'].fillna(credit_card_data['Is Fraud?'].mode()[0], inplace=True)
```

As discussed earlier:

- The missing values in the 'Merchant State' column were filled with the placeholder 'Unknown', ensuring that all entries had a valid state identifier.
- For the 'Zip' column, a custom approach was applied. A dictionary mapping, named 'element\_dict', was created to associate missing zip codes with their respective 'Merchant City' entries. This was done using a lambda function applied row-wise to the dataframe. If a zip code was missing (pd.isna(row['Zip'])), it was replaced with the corresponding value from 'element\_dict' based on the 'Merchant City' entry. This method allowed for precise imputation of missing zip codes based on the known city information.
- Missing values in the 'Errors?' column were replaced with 'No Error', assuming that a

```
In [73]: 1 credit_card_data['Merchant State'].fillna('Unknown', inplace=True)
2 credit_card_data['Zip'] = credit_card_data.apply(lambda row: element_dict[row['Merchant City']], axis=1)
3 credit_card_data['Errors?'].fillna('No Error', inplace=True)
```

## Converting 'Time' to minutes since midnight.

```
In [75]: 1 # Function to convert time to minutes since midnight
2 def convert_time(time_str):
3     try:
4         t = pd.to_datetime(time_str, format='%H:%M')
5         return t.hour * 60 + t.minute
6     except:
7         return None
8     # Returns None for any incorrectly formatted time strings
9
10    # Apply the conversion function to the 'Time' column
11    credit_card_data['Time'] = credit_card_data['Time'].apply(convert_time)
```

## Preprocessing Results

In [74]: 1 credit\_card\_data.isnull().sum()

Out[74]:

User	0
Card	0
Year	0
Month	0
Day	0
Time	0
Amount	0
Use Chip	0
Merchant Name	0
Merchant City	0
Merchant State	0
Zip	0
MCC	0
Errors?	0
Is Fraud?	0
dtype: int64	

In [76]: 1 credit\_card\_data

Out[76]:

	User	Card	Year	Month	Day	Time	Amount	Use Chip	Merchant Name	Merchant City
0	0	0	2002	9	1	381	\$134.09	Swipe Transaction	3.527213e+18	La Verne
1	0	0	2002	9	1	402	\$38.48	Swipe Transaction	-7.276121e+17	Monterey Park
2	0	0	2002	9	2	382	\$120.34	Swipe Transaction	-7.276121e+17	Monterey Park
3	0	0	2002	9	2	1065	\$128.95	Swipe Transaction	3.414527e+18	Monterey Park
4	0	0	2002	9	3	383	\$104.71	Swipe Transaction	5.817218e+18	La Verne
...	...	...	...	...	...	...	...	...	...	...
2555185	215	0	2006	11	14	439	\$55.43	Swipe Transaction	4.591017e+18	Whitehouse Station
2555186	215	0	2006	11	14	472	\$7.23	Swipe Transaction	5.205100e+18	North Brunswick
2555187	215	0	2006	11	14	535	\$79.07	Swipe Transaction	8.384250e+17	New York
2555188	215	0	2006	11	14	667	\$22.05	Swipe Transaction	5.205100e+18	North Brunswick
2555189	215	0	2006	11	14	1	\$80.00	Swipe Transaction	1.799190e+18	ONLINE

2555190 rows × 15 columns



## Data Preprocessing and Encoding Categorical Variables

- This code segment performs data preprocessing and encoding of categorical variables in the credit card transaction dataset. The 'Amount' column is converted to numeric format by removing the dollar sign. Categorical columns ('Use Chip', 'Merchant City', 'Merchant State', 'Errors?', 'Is Fraud?') are label encoded using scikit-learn's LabelEncoder. Additionally, the 'Merchant Name' column is converted from scientific notation to an integer format if necessary and then label encoded. Finally, the preprocessed data is saved to a new CSV file named 'card\_transaction\_processed.csv'.

In [77]:

```

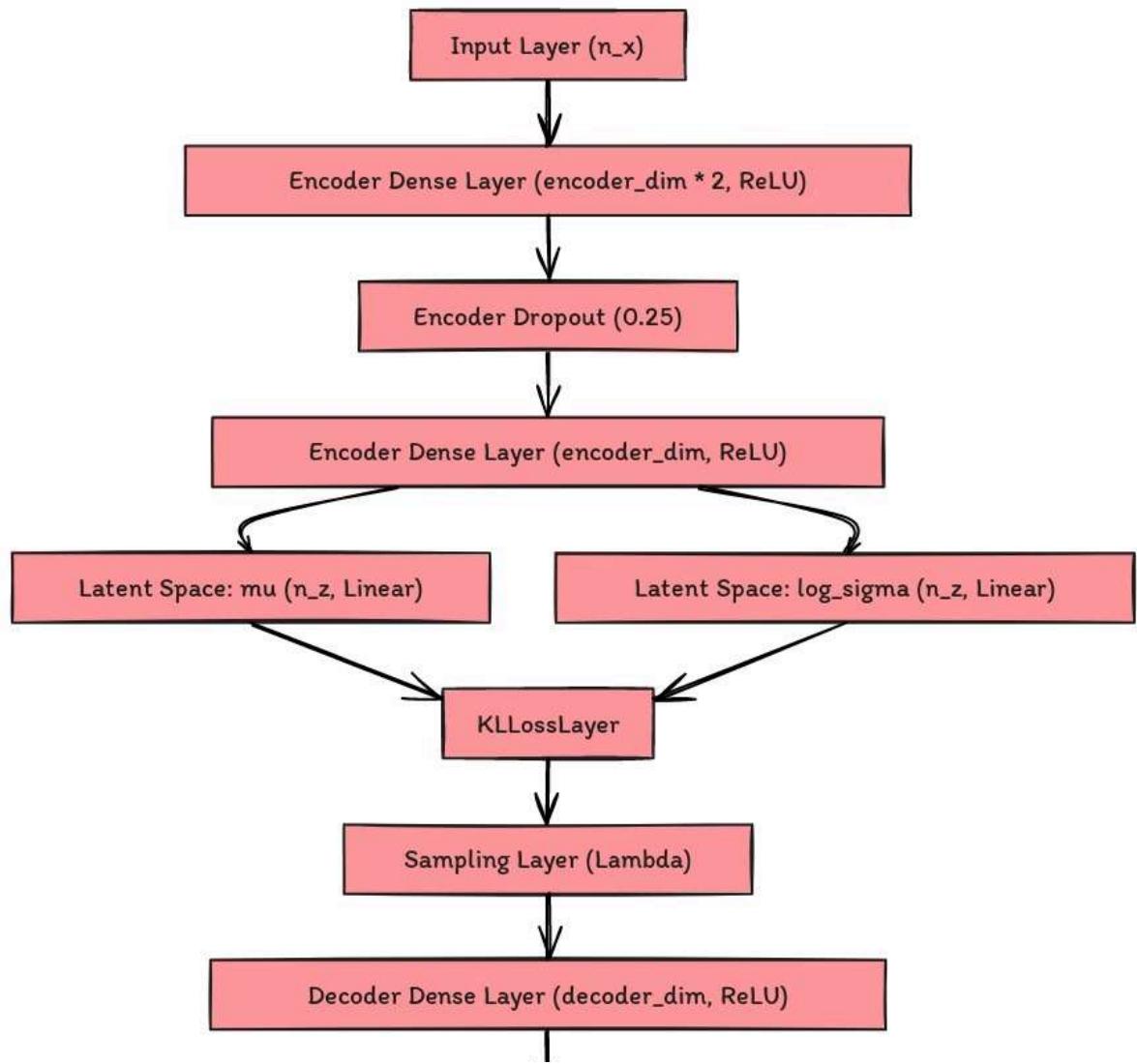
1 # Convert 'Amount' to numeric after stripping the dollar sign
2 credit_card_data['Amount'] = credit_card_data['Amount'].replace('[\$,]', ''
3
4 # Label encode categorical columns
5
6 categorical_columns = ['Use Chip', 'Merchant City', 'Merchant State', 'Er
7 encoders = {col: LabelEncoder() for col in categorical_columns}
8 for column in categorical_columns:
9     credit_card_data[column] = encoders[column].fit_transform(credit_card_
10
11 # Convert 'Merchant Name' from scientific notation to integer if needed ar
12 # This depends on whether 'Merchant Name' is a floating point representati
13 # Here is a generic approach if 'Merchant Name' is indeed a large floating
14 encoder['Merchant Name'] = LabelEncoder()
15 credit_card_data['Merchant Name'] = credit_card_data['Merchant Name'].apply(
16 credit_card_data['Merchant Name'] = label_encoder.fit_transform(credit_car
17
18 # Save the preprocessed data to a new CSV file
19 credit_card_data.to_csv('card_transaction_processed.csv', index=False)
20
21 print("Data preprocessing complete. The processed data is saved to 'card_t

```

Data preprocessing complete. The processed data is saved to 'card\_transaction\_processed.csv'.

## Varational Autoencoder for Finance-Tabular synthetic data generation

- The architecture of the VAE that we have used is visualized below:



## Normalizing data for a good life

- This code loads data, splits it into training/validation sets, and standardizes numeric columns for machine learning preprocessing.

In [79]:

```
1 # Load the data
2 data = pd.read_csv('card_transaction_processed.csv')
3
4 # Proceed with your preprocessing...
5 categorical_columns = ['Use Chip', 'Merchant City', 'Merchant State', 'Zip'
6 numeric_cols = ['Time', 'Amount', 'Merchant Name']
7
8 X = data
9
10 # Split data into training and validation sets
11 X_train, X_valid = train_test_split(X, test_size=0.2, random_state=42)
12
13 # Initialize and apply scalers
14 scaler = StandardScaler()
15 X_train_scaled = X_train.copy()
16 X_valid_scaled = X_valid.copy()
17 X_train_scaled = pd.DataFrame(scaler.fit_transform(X_train), columns=X_train.co
18 X_valid_scaled = pd.DataFrame(scaler.transform(X_valid), columns=X_train.co
```

```
In [80]: 1 # Checking for NaN values in the dataset
2 print(X_train_scaled.isnull().sum())
3 print('-----')
4 print(X_valid_scaled.isnull().sum())
```

```
User          0
Card          0
Year          0
Month         0
Day           0
Time          0
Amount         0
Use Chip      0
Merchant Name 0
Merchant City 0
Merchant State 0
Zip           0
MCC            0
Errors?        0
Is Fraud?     0
dtype: int64
-----
User          0
Card          0
Year          0
Month         0
Day           0
Time          0
Amount         0
Use Chip      0
Merchant Name 0
Merchant City 0
Merchant State 0
Zip           0
MCC            0
Errors?        0
Is Fraud?     0
dtype: int64
```

## The CRUX!!

- This code defines a Variational Autoencoder (VAE) with custom layers for KL loss computation, utilizing a more complex architecture with dropout regularization and advanced decoder layers. The VAE loss combines reconstruction loss with the KL divergence added by the custom layer, and the model is compiled with an Adam optimizer and specified learning rate for training.

### The objection Function:

$$\text{VAE Loss} = \text{Reconstruction Loss} + \text{KL Divergence Loss}$$

$$\text{KL Divergence Loss} = -0.5 \sum_{j=1}^{n_z} (1 + \log(\sigma_j^2) - \mu_j^2 - \sigma_j^2)$$

$$\text{Reconstruction Loss} = \frac{1}{N} \sum_{i=1}^N (y_{\text{true}}^{(i)} - y_{\text{pred}}^{(i)})^2$$

*Where :*

$N$  be the number of samples in the batch.

$y_{\text{true}}(i)$  be the true output for the  $i$ -th sample.

$y_{\text{pred}}(i)$  be the predicted output for the  $i$ -th sample.

In [81]:

```

1 # Hyperparameters
2 m = 1000 # batch size
3 n_z = 2 # Latent space size
4 encoder_dim = 5 # dim of encoder hidden layer
5 decoder_dim = 5 # dim of decoder hidden layer
6 n_x = X_train_scaled.shape[1] # Input dimension
7 n_y = 1 # Output dimension for Label (fraud/not fraud as additional input)
8
9 from tensorflow.keras.layers import Layer
10
11 class KLossLayer(Layer):
12     """ Custom layer to compute KL loss """
13     def __init__(self, **kwargs):
14         super(KLossLayer, self).__init__(**kwargs)
15
16     def call(self, inputs):
17         mu, log_sigma = inputs
18         kl_loss = -0.5 * K.sum(1 + log_sigma - K.square(mu) - K.exp(log_sigma))
19         self.add_loss(K.mean(kl_loss)) # Add KL divergence to the model's loss
20         return inputs
21
22
23 # Redefine VAE with a more complex architecture
24 inputs = Input(shape=(n_x,))
25 encoder_h = Dense(encoder_dim*2, activation='relu')(inputs)
26 encoder_h = Dropout(0.25)(encoder_h) # Add dropout for regularization
27 encoder_h = Dense(encoder_dim, activation='relu')(encoder_h) # Increase dimension
28 mu = Dense(n_z, activation='linear')(encoder_h)
29 log_sigma = Dense(n_z, activation='linear')(encoder_h)
30
31 # Sampling function
32 def sample_z(args):
33     mu, log_sigma = args
34     eps = K.random_normal(shape=(K.shape(mu)[0], n_z), mean=0., stddev=1.)
35     return mu + K.exp(log_sigma / 2) * eps
36
37 # Use Lambda Layer for sampling
38 z = Lambda(sample_z, output_shape=(n_z,))([mu, log_sigma])
39
40 # KL Loss layer added
41 mu, log_sigma = KLossLayer()([mu, log_sigma])
42
43 # More complex decoder
44 decoder_h = Dense(decoder_dim, activation='relu')(z)
45 decoder_h = Dropout(0.25)(decoder_h) # Add dropout for regularization
46 decoder_h = Dense(decoder_dim*2, activation='relu')(decoder_h)
47 y_pred = Dense(n_x, activation='sigmoid')(decoder_h)
48
49 # VAE model
50 vae = Model(inputs, y_pred)
51
52 # Define the VAE Loss (only reconstruction Loss here since KL Loss is added)
53 def vae_loss(y_true, y_pred, alpha = 0.1):
54     reconstruction_loss = K.mean(K.square(y_true - y_pred))
55     # reconstruction_loss = K.mean(categorical_crossentropy(y_true, y_pred))
56     return reconstruction_loss
57

```

```
58  
59 # Example usage:  
60  
61 # Compile the VAE  
62 # Adjust learning rate  
63 vae.compile(optimizer=Adam(learning_rate=0.0001), loss=vae_loss)  
64
```

In [82]: 1 vae.summary()

Model: "model\_2"

Layer (type)	Output Shape	Param #	Connected to
<hr/>			
<hr/>			
input_3 (InputLayer)	[(None, 15)]	0	[]
dense_7 (Dense)	(None, 10)	160	['input_3[0][0]']
dropout_2 (Dropout)	(None, 10)	0	['dense_7[0][0]']
dense_8 (Dense)	(None, 5)	55	['dropout_2[0][0]']
dense_9 (Dense)	(None, 2)	12	['dense_8[0][0]']
dense_10 (Dense)	(None, 2)	12	['dense_8[0][0]']
lambda_1 (Lambda)	(None, 2)	0	['dense_9[0][0]', 'dense_10[0][0]']
dense_11 (Dense)	(None, 5)	15	['lambda_1[0][0]']
dropout_3 (Dropout)	(None, 5)	0	['dense_11[0][0]']
dense_12 (Dense)	(None, 10)	60	['dropout_3[0][0]']
dense_13 (Dense)	(None, 15)	165	['dense_12[0][0]']
<hr/>			
<hr/>			
Total params: 479 (1.87 KB)			
Trainable params: 479 (1.87 KB)			
Non-trainable params: 0 (0.00 Byte)			

## Training

- The code creates a learning rate scheduler in TensorFlow for dynamic adjustment of the learning rate during training. It utilizes an exponential decay strategy after an initial phase

In [83]:

```
1 import tensorflow as tf
2
3 def scheduler(epoch, lr):
4     if epoch < 10:
5         return lr
6     else:
7         # Compute the new Learning rate
8         new_lr = lr * tf.math.exp(-0.1).numpy() # Convert tensor to numpy
9         return float(new_lr) # Ensure the return type is float
10
11 # Create the Learning rate scheduler callback
12 lr_scheduler = LearningRateScheduler(scheduler)
13
14 # Now, use this callback in your model training
15 history = vae.fit(X_train_scaled, X_train_scaled,
16                     batch_size=m,
17                     epochs=50,
18                     verbose=1,
19                     validation_data=(X_valid_scaled, X_valid_scaled),
20                     callbacks=[EarlyStopping(patience=5), lr_scheduler])
```

```
Epoch 1/50
2045/2045 [=====] - 15s 6ms/step - loss: 1.1084 - va
l_loss: 0.9921 - lr: 1.0000e-04
Epoch 2/50
2045/2045 [=====] - 12s 6ms/step - loss: 0.9920 - va
l_loss: 0.9630 - lr: 1.0000e-04
Epoch 3/50
2045/2045 [=====] - 11s 6ms/step - loss: 0.9694 - va
l_loss: 0.9428 - lr: 1.0000e-04
Epoch 4/50
2045/2045 [=====] - 10s 5ms/step - loss: 0.9557 - va
l_loss: 0.9328 - lr: 1.0000e-04
Epoch 5/50
2045/2045 [=====] - 12s 6ms/step - loss: 0.9475 - va
l_loss: 0.9256 - lr: 1.0000e-04
Epoch 6/50
2045/2045 [=====] - 12s 6ms/step - loss: 0.9421 - va
l_loss: 0.9210 - lr: 1.0000e-04
Epoch 7/50
2045/2045 [=====] - 12s 6ms/step - loss: 0.9387 - va
l_loss: 0.9180 - lr: 1.0000e-04
Epoch 8/50
2045/2045 [=====] - 12s 6ms/step - loss: 0.9362 - va
l_loss: 0.9155 - lr: 1.0000e-04
Epoch 9/50
2045/2045 [=====] - 10s 5ms/step - loss: 0.9340 - va
l_loss: 0.9130 - lr: 1.0000e-04
Epoch 10/50
2045/2045 [=====] - 10s 5ms/step - loss: 0.9319 - va
l_loss: 0.9098 - lr: 1.0000e-04
Epoch 11/50
2045/2045 [=====] - 12s 6ms/step - loss: 0.9296 - va
l_loss: 0.9064 - lr: 9.0484e-05
Epoch 12/50
2045/2045 [=====] - 11s 6ms/step - loss: 0.9274 - va
l_loss: 0.9041 - lr: 8.1873e-05
Epoch 13/50
2045/2045 [=====] - 12s 6ms/step - loss: 0.9257 - va
l_loss: 0.9023 - lr: 7.4082e-05
Epoch 14/50
2045/2045 [=====] - 12s 6ms/step - loss: 0.9244 - va
l_loss: 0.9007 - lr: 6.7032e-05
Epoch 15/50
2045/2045 [=====] - 11s 6ms/step - loss: 0.9234 - va
l_loss: 0.8998 - lr: 6.0653e-05
Epoch 16/50
2045/2045 [=====] - 12s 6ms/step - loss: 0.9226 - va
l_loss: 0.8988 - lr: 5.4881e-05
Epoch 17/50
2045/2045 [=====] - 15s 8ms/step - loss: 0.9221 - va
l_loss: 0.8980 - lr: 4.9659e-05
Epoch 18/50
2045/2045 [=====] - 12s 6ms/step - loss: 0.9216 - va
l_loss: 0.8974 - lr: 4.4933e-05
Epoch 19/50
2045/2045 [=====] - 12s 6ms/step - loss: 0.9212 - va
l_loss: 0.8969 - lr: 4.0657e-05
```

```
Epoch 20/50
2045/2045 [=====] - 12s 6ms/step - loss: 0.9208 - va
l_loss: 0.8963 - lr: 3.6788e-05
Epoch 21/50
2045/2045 [=====] - 12s 6ms/step - loss: 0.9205 - va
l_loss: 0.8958 - lr: 3.3287e-05
Epoch 22/50
2045/2045 [=====] - 12s 6ms/step - loss: 0.9202 - va
l_loss: 0.8955 - lr: 3.0119e-05
Epoch 23/50
2045/2045 [=====] - 12s 6ms/step - loss: 0.9200 - va
l_loss: 0.8952 - lr: 2.7253e-05
Epoch 24/50
2045/2045 [=====] - 11s 6ms/step - loss: 0.9199 - va
l_loss: 0.8950 - lr: 2.4660e-05
Epoch 25/50
2045/2045 [=====] - 11s 6ms/step - loss: 0.9195 - va
l_loss: 0.8947 - lr: 2.2313e-05
Epoch 26/50
2045/2045 [=====] - 12s 6ms/step - loss: 0.9194 - va
l_loss: 0.8945 - lr: 2.0190e-05
Epoch 27/50
2045/2045 [=====] - 11s 6ms/step - loss: 0.9193 - va
l_loss: 0.8942 - lr: 1.8268e-05
Epoch 28/50
2045/2045 [=====] - 12s 6ms/step - loss: 0.9191 - va
l_loss: 0.8939 - lr: 1.6530e-05
Epoch 29/50
2045/2045 [=====] - 12s 6ms/step - loss: 0.9190 - va
l_loss: 0.8939 - lr: 1.4957e-05
Epoch 30/50
2045/2045 [=====] - 11s 6ms/step - loss: 0.9188 - va
l_loss: 0.8937 - lr: 1.3534e-05
Epoch 31/50
2045/2045 [=====] - 11s 5ms/step - loss: 0.9188 - va
l_loss: 0.8935 - lr: 1.2246e-05
Epoch 32/50
2045/2045 [=====] - 11s 6ms/step - loss: 0.9188 - va
l_loss: 0.8934 - lr: 1.1080e-05
Epoch 33/50
2045/2045 [=====] - 11s 5ms/step - loss: 0.9186 - va
l_loss: 0.8934 - lr: 1.0026e-05
Epoch 34/50
2045/2045 [=====] - 13s 7ms/step - loss: 0.9185 - va
l_loss: 0.8933 - lr: 9.0718e-06
Epoch 35/50
2045/2045 [=====] - 12s 6ms/step - loss: 0.9185 - va
l_loss: 0.8932 - lr: 8.2085e-06
Epoch 36/50
2045/2045 [=====] - 11s 5ms/step - loss: 0.9184 - va
l_loss: 0.8931 - lr: 7.4274e-06
Epoch 37/50
2045/2045 [=====] - 11s 5ms/step - loss: 0.9183 - va
l_loss: 0.8931 - lr: 6.7206e-06
Epoch 38/50
2045/2045 [=====] - 11s 5ms/step - loss: 0.9183 - va
l_loss: 0.8930 - lr: 6.0810e-06
```

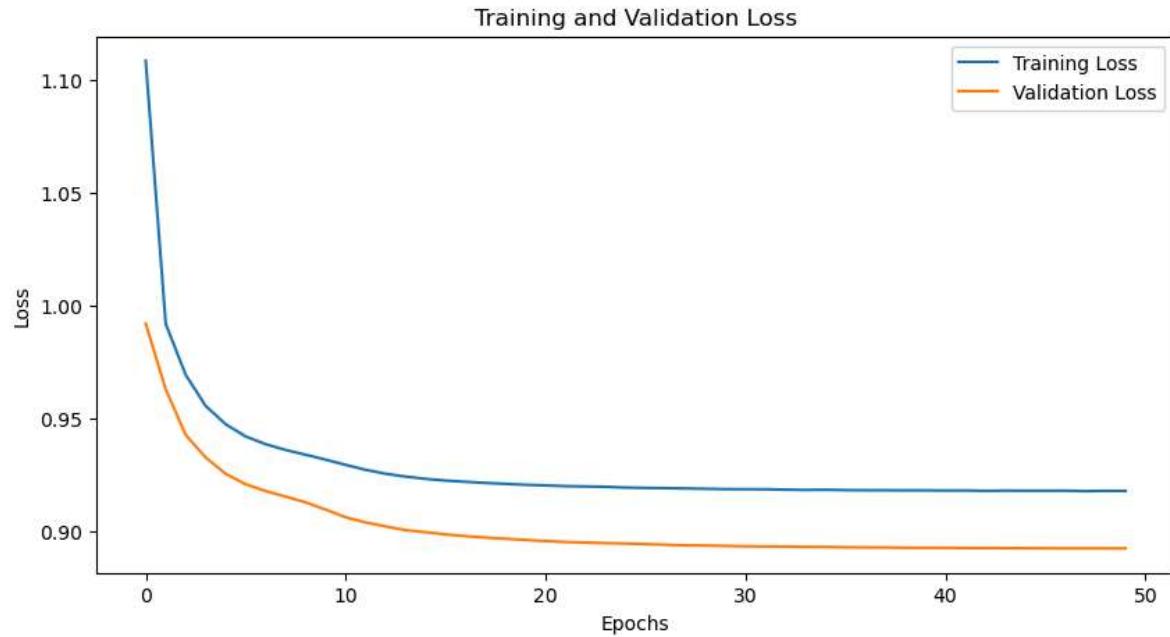
```
Epoch 39/50
2045/2045 [=====] - 11s 6ms/step - loss: 0.9183 - va
l_loss: 0.8929 - lr: 5.5023e-06
Epoch 40/50
2045/2045 [=====] - 11s 5ms/step - loss: 0.9183 - va
l_loss: 0.8929 - lr: 4.9787e-06
Epoch 41/50
2045/2045 [=====] - 11s 5ms/step - loss: 0.9182 - va
l_loss: 0.8928 - lr: 4.5049e-06
Epoch 42/50
2045/2045 [=====] - 11s 6ms/step - loss: 0.9182 - va
l_loss: 0.8928 - lr: 4.0762e-06
Epoch 43/50
2045/2045 [=====] - 12s 6ms/step - loss: 0.9181 - va
l_loss: 0.8928 - lr: 3.6883e-06
Epoch 44/50
2045/2045 [=====] - 14s 7ms/step - loss: 0.9182 - va
l_loss: 0.8927 - lr: 3.3373e-06
Epoch 45/50
2045/2045 [=====] - 11s 5ms/step - loss: 0.9181 - va
l_loss: 0.8927 - lr: 3.0197e-06
Epoch 46/50
2045/2045 [=====] - 11s 5ms/step - loss: 0.9181 - va
l_loss: 0.8927 - lr: 2.7324e-06
Epoch 47/50
2045/2045 [=====] - 11s 5ms/step - loss: 0.9181 - va
l_loss: 0.8926 - lr: 2.4724e-06
Epoch 48/50
2045/2045 [=====] - 11s 5ms/step - loss: 0.9180 - va
l_loss: 0.8926 - lr: 2.2371e-06
Epoch 49/50
2045/2045 [=====] - 12s 6ms/step - loss: 0.9181 - va
l_loss: 0.8926 - lr: 2.0242e-06
Epoch 50/50
2045/2045 [=====] - 11s 5ms/step - loss: 0.9180 - va
l_loss: 0.8926 - lr: 1.8316e-06
```

## Training v/s Validation Loss

- This code snippet visualizes the training and validation loss over epochs, helping to analyze the model's performance and convergence during training.

In [84]:

```
1 import matplotlib.pyplot as plt
2
3 # Plotting the training and validation loss
4 plt.figure(figsize=(10, 5))
5 plt.plot(history.history['loss'], label='Training Loss')
6 plt.plot(history.history['val_loss'], label='Validation Loss')
7 plt.title('Training and Validation Loss')
8 plt.xlabel('Epochs')
9 plt.ylabel('Loss')
10 plt.legend()
11 plt.show()
```



## Encoder Model Creation and Data Transformation

- The code defines an encoder model using layers from an existing Variational Autoencoder (VAE) model. It extracts the latent representations of the input data by passing it through the encoder layers. The encoder\_model is then used to transform the training data into the latent space, which is crucial for tasks like dimensionality reduction or feature extraction.

```
In [ ]: 1 ### Create a new input Layer that matches the input dimension of the VAE
2 encoder_input = Input(shape=(n_x,), name='encoder_input')
3
4 # Reuse layers from the existing VAE model to define the encoder
5 x = encoder_input
6
7 # Apply each appropriate layer of the VAE to build up to the Lambda_4 Layer
8 for layer in vae.layers[1:4]: # This includes input_layer to dense_30
9     x = layer(x)
10
11 # Now manually connect dense_31 and dense_32 from the output of dense_30
12 mu = vae.layers[4](x)
13 log_sigma = vae.layers[5](x)
14
15 # The Lambda Layer takes both mu and Log_sigma as inputs
16 z = vae.layers[6]([mu, log_sigma])
17
18 # Create the encoder model that outputs the result from the Lambda_4 Layer
19 encoder_model = Model(inputs=encoder_input, outputs=x, name='encoder_model')
20
21 # Display the model summary to verify the architecture
22 encoder_model.summary()
23
24
25 # Now use this model to transform your data into the Latent space
26 latent_representations = encoder_model.predict(X_train_scaled)
27
```

## Latent Space Visualization

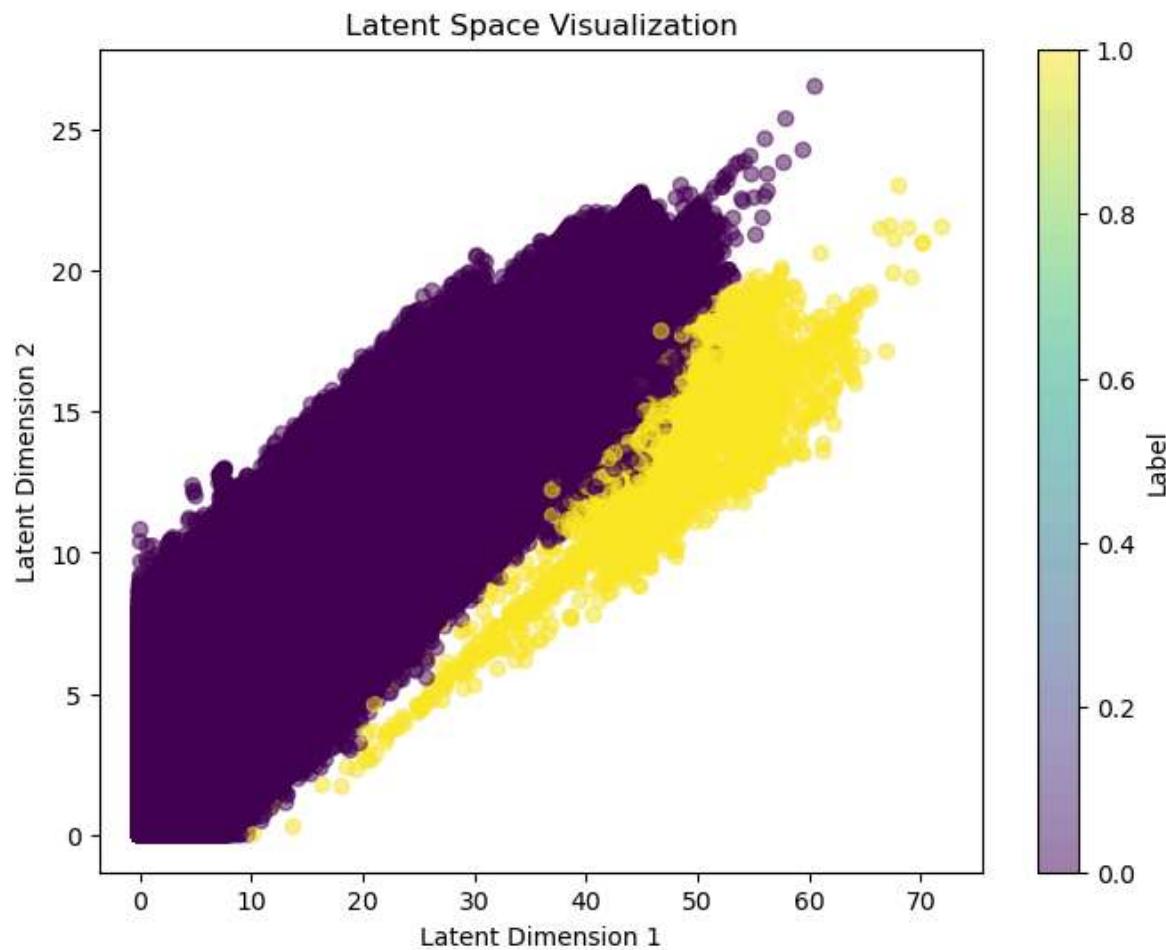
- The code snippet creates a scatter plot to visualize the latent representations of data points in a 2-dimensional latent space. The points are colored based on their true category labels from y\_train using same seed as used earlier for train/valid split. This visualization helps in understanding the clustering or distribution of data points in the latent space, which is essential for interpreting the learned representations in a VAE.

In [86]:

```

1 import matplotlib.pyplot as plt
2
3 # Assuming you have some labels in y_train to color the points according to
4 # If you don't have such labels, you can simply plot without coloring
5 plt.figure(figsize=(8, 6))
6 # Separate features and target
7 X = data.drop('Is Fraud?', axis=1)
8 y = data['Is Fraud?']
9 X_train, X_valid, y_train, y_valid = train_test_split(X, y, test_size=0.2)
10 plt.scatter(latent_representations[:, 0], latent_representations[:, 1], c=y)
11 plt.colorbar(label='Label')
12 plt.xlabel('Latent Dimension 1')
13 plt.ylabel('Latent Dimension 2')
14 plt.title('Latent Space Visualization')
15 plt.show()

```



## Synthetic Data Generation using VAE Decoder

- The code snippet defines a decoder model using the trained VAE model to generate synthetic data. It creates synthetic data points by sampling from the latent space and passing them through the decoder layers. The `generate_synthetic_data` function allows adjusting the temperature parameter for controlling the diversity of the generated data. This approach enables the generation of new data points similar to the original dataset, useful for augmenting data or generating synthetic examples for testing models.

## Justification for Synthetic Data Generation

Synthetic data generation using a Variational Autoencoder (VAE) decoder is a valuable technique in data science for several reasons:

1. **Data Augmentation:** It allows for increasing the size of the dataset by generating new data points that follow the same distribution as the original data. This can be particularly useful when dealing with limited data availability.
2. **Diversity Exploration:** By adjusting the temperature parameter, the diversity of the synthetic data can be controlled. Higher temperatures result in more diverse data, enabling exploration of different scenarios and edge cases.
3. **Model Testing:** Synthetic data can be used for testing machine learning models, especially in cases where obtaining real-world data is challenging or costly. It helps in evaluating the model's generalization and robustness.
4. **Privacy Preservation:** Synthetic data generation can also be used for privacy-preserving data sharing. It allows sharing insights or models without exposing sensitive or personally identifiable information present in the original dataset.

Overall, leveraging VAE-based synthetic data generation enhances data-driven workflows by providing additional data points for training, testing, and exploration while maintaining the underlying data characteristics.

In [87]:

```

1  from tensorflow.keras.models import Model
2
3  # Assuming `vae` is your trained model and you know the Layer where the de
4  # For instance, if your decoder starts at Layer 10:
5  decoder_input = Input(shape=(n_z,))
6  x = decoder_input
7
8  # Reapply each decoder layer
9  for layer in vae.layers[7:]: # adjust the index 10 to where your decoder
10    x = layer(x)
11
12 decoder = Model(decoder_input, x)
13
14 # Now you can use this decoder to generate synthetic data
15 def generate_synthetic_data(decoder, num_samples, temperature=1.0):
16     z_sample = np.random.normal(scale=temperature, size=(num_samples, n_z))
17     synthetic_data = decoder.predict(z_sample)
18     return synthetic_data
19
20 # Generate data with a higher temperature for more diversity
21 temperature = 0.75 # Lowering temperature might tighten the distribution
22 num_samples = 2000000 # Number of synthetic transactions to generate
23 synthetic_data = generate_synthetic_data(decoder, num_samples, temperature)

```

62500/62500 [=====] - 107s 2ms/step

```
In [88]: 1 df_gen = pd.DataFrame(synthetic_data, columns=data.columns)
2 df_gen.head()
```

Out[88]:

	User	Card	Year	Month	Day	Time	Amount	Use Chip	Merchant I Name
0	0.000067	0.995171	2.048981e-07	0.012143	0.002433	0.651233	0.073876	0.196313	0.020216
1	0.000278	0.988760	2.965598e-06	0.012659	0.003737	0.564943	0.059543	0.309353	0.062197
2	0.000003	0.979297	4.345837e-09	0.003396	0.067368	0.045733	0.006237	0.268728	0.072838
3	0.000009	0.998046	9.809742e-09	0.007216	0.001507	0.649788	0.054341	0.125381	0.013765
4	0.000189	0.992039	1.022046e-06	0.015412	0.003290	0.654292	0.085520	0.243257	0.026869

## Denormalization of Data

- The denormalization process involves reverting the scaled or normalized data back to its original scale, which is essential for interpreting and using the data correctly. In this code snippet, the `scaler.inverse_transform` function is applied to `df_gen`, which is assumed to be the original DataFrame before scaling. The resulting `denormalizedDf` DataFrame contains the denormalized data.

```
In [89]: 1 denormalizedDf = scaler.inverse_transform(df_gen)
```

```
2
3 # Assuming df_gen is your original DataFrame and 'denormalizedDf' is the result
4 denormalizedDf = pd.DataFrame(denormalizedDf, columns=df_gen.columns)
5
6 # Now you can use .head() on this DataFrame
7 denormalizedDf.head()
```

Out[89]:

	User	Card	Year	Month	Day	Time	Amount	Use Chip	
0	102.066132	2.577117	2011.792358	6.562517	15.738690	974.997498	50.573677	1.544594	103
1	102.079254	2.568495	2011.792358	6.564306	15.750158	948.481506	49.371899	1.641368	106
2	102.062134	2.555768	2011.792358	6.532155	16.309546	788.932922	44.902348	1.606589	106
3	102.062531	2.580983	2011.792358	6.545414	15.730554	974.553528	48.935707	1.483868	103
4	102.073738	2.572905	2011.792358	6.573864	15.746228	975.937500	51.549957	1.584783	103

## Rounding Values in DataFrame

- The code snippet uses `applymap` along with a lambda function to round each value in the `denormalizedDf` DataFrame. If the decimal part of a value is greater than or equal to 0.5,

it is rounded up using `np.ceil`, otherwise it is rounded down using `np.floor`. This operation ensures that the data is rounded to the nearest integer for each value in the DataFrame.

```
In [90]: 1 df = denormalizedDf.applymap(lambda x: np.ceil(x) if x - np.floor(x) >= 0
```

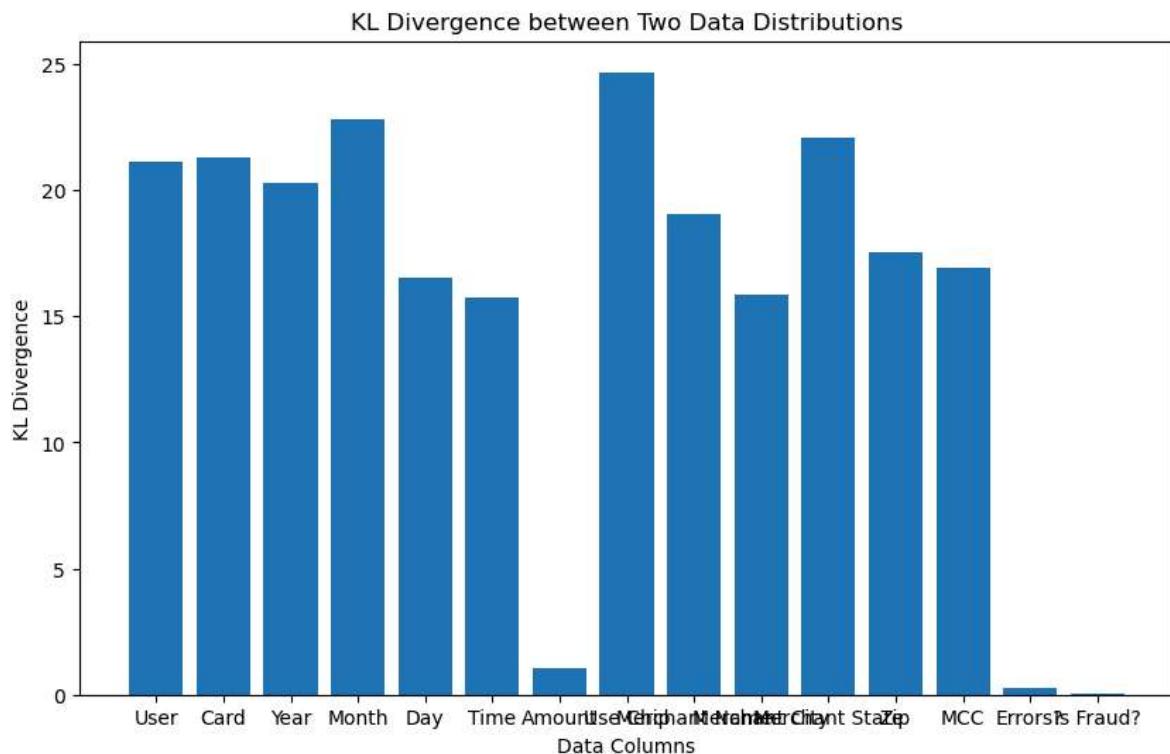
## KL Divergence Analysis

- The provided code defines a function `k1_divergence` to calculate the Kullback-Leibler (KL) divergence between two probability distributions. The `calculate_kl_divergence` function computes the KL divergence for each numerical column in two DataFrames `df1` and `df2`, assuming they are preprocessed and scaled. The resulting KL divergences are then visualized using a bar chart, showcasing the divergence between the distributions of the two datasets across different columns.

In [92]:

```
1 import pandas as pd
2 import numpy as np
3 import matplotlib.pyplot as plt
4
5 # Define the KL divergence function
6 def kl_divergence(p, q, eps=1e-10):
7     p = np.array(p) + eps # Add eps to p
8     q = np.array(q) + eps # Add eps to q
9
10    p = p / np.sum(p)
11    q = q / np.sum(q)
12
13    return np.sum(p * np.log(p / q))
14
15
16 # Calculate KL divergence for each column in the DataFrames
17 def calculate_kl_divergence(df1, df2, bins=30):
18     kl_divergences = {}
19     for column in df1.columns:
20         if df1[column].dtype != 'object' and df2[column].dtype != 'object':
21             hist1, bin_edges = np.histogram(df1[column], bins=bins, density=True)
22             hist2, _ = np.histogram(df2[column], bins=bin_edges, density=True)
23             kl_div = kl_divergence(hist1, hist2)
24             kl_divergences[column] = kl_div
25     return kl_divergences
26
27 X = pd.DataFrame(scaler.fit_transform(data), columns=data.columns)
28
29 # Calculate KL divergences
30 kl_divergences = calculate_kl_divergence(X, df_gen)
31
32 # Plotting the KL divergences
33 plt.figure(figsize=(10, 6))
34 plt.bar(kl_divergences.keys(), kl_divergences.values())
35 plt.xlabel('Data Columns')
36 plt.ylabel('KL Divergence')
37 plt.title('KL Divergence between Two Data Distributions')
38 plt.show()
```





## Decoding Categorical Data

- The provided code snippet decodes categorical data columns using inverse transformations from previously defined encoders. It iterates through each column in the DataFrame `data`, checks if the column is encoded, and performs the inverse transformation to retrieve the original labels. The resulting DataFrame `denormalizedDf_Labeled` contains the decoded categorical columns along with the unchanged numerical columns.

```
In [116]: 1 denormalizedDf_Labeled = pd.DataFrame()
2 for col in data.columns:
3     try:
4         if col in encoders:
5             denormalizedDf_Labeled[col] = encoders[col].inverse_transform(
6                 data[[col]])
7         else:
8             denormalizedDf_Labeled[col] = df[col]
9     except Exception as e:
10         print(f"Error decoding {col}: {str(e)}")
```

```
In [117]: 1 # Convert numerical columns from float to int
2 denormalizedDf_Labeled['User'] = denormalizedDf_Labeled['User'].astype(int)
3 denormalizedDf_Labeled['Card'] = denormalizedDf_Labeled['Card'].astype(int)
4 denormalizedDf_Labeled['Year'] = denormalizedDf_Labeled['Year'].astype(int)
5 denormalizedDf_Labeled['Month'] = denormalizedDf_Labeled['Month'].astype(int)
6 denormalizedDf_Labeled['Day'] = denormalizedDf_Labeled['Day'].astype(int)
7 denormalizedDf_Labeled['Time'] = denormalizedDf_Labeled['Time'].astype(int)
8 denormalizedDf_Labeled['Amount'] = denormalizedDf_Labeled['Amount'].astype(float)
9 denormalizedDf_Labeled['Merchant Name'] = denormalizedDf_Labeled['Merchant Name'].str.replace(' ', '_')
10 denormalizedDf_Labeled['Zip'] = denormalizedDf_Labeled['Zip'].astype(int)
11 denormalizedDf_Labeled['MCC'] = denormalizedDf_Labeled['MCC'].astype(int)
```

```
In [119]: 1 # Convert 'Time' column to HH:MM format
2 denormalizedDf_Labeled['Time'] = pd.to_datetime(denormalizedDf_Labeled['Time']).dt.strftime('%H:%M')
```

```
In [121]: 1 # Add '$' before every value in the 'Amount' column
2 denormalizedDf_Labeled['Amount'] = '$' + denormalizedDf_Labeled['Amount'].map(str)
```

## Results

```
In [122]: 1 # Display the first few rows to verify decoding
2 denormalizedDf_Labeled.head()
```

Out[122]:

	User	Card	Year	Month	Day	Time	Amount	Use Chip	Merchant Name	Merchant City	Merchant State
0	102	3	2012	7	16	16:15	\$51	Swipe Transaction	10346	Melville	Macedonia
1	102	3	2012	7	16	15:48	\$49	Swipe Transaction	10614	Mountain Grove	Macedonia
2	102	3	2012	7	16	13:09	\$45	Swipe Transaction	10682	Ocoee	Macedonia
3	102	3	2012	7	16	16:15	\$49	Online Transaction	10304	Marlboro	Macedonia
4	102	3	2012	7	16	16:16	\$52	Swipe Transaction	10388	Milladore	Macedonia

### 1. Coarse grained

- The metric help gauge the level of similarity and uniqueness between the real and synthetic datasets, providing insights into the quality and diversity of the synthetic data generation process.
- It uses two foundations for evalution:

#### 1. Direct Copy Percentage

- It calculates what percentage of the combined (real and synthetic) data is a direct copy of the real data.
- Direct Copy Percentage =  $\frac{\text{Total Data Size}}{\text{Number of Copies}} \times 100$

- Here, "copies" refer to rows in the combined data that are not present in either the real or synthetic data.

## 2. Self-Copy Percentage:

- It measures the percentage of duplicate rows (self-copies) within the synthetic data.
- Self-Copy Percentage =  $\frac{\text{Number of Duplicate Rows in Synthetic Data}}{\text{Total Synthetic Data Size}} \times 100$
- This metric reveals how much of the synthetic data consists of repeated patterns or copies of existing rows within itself.

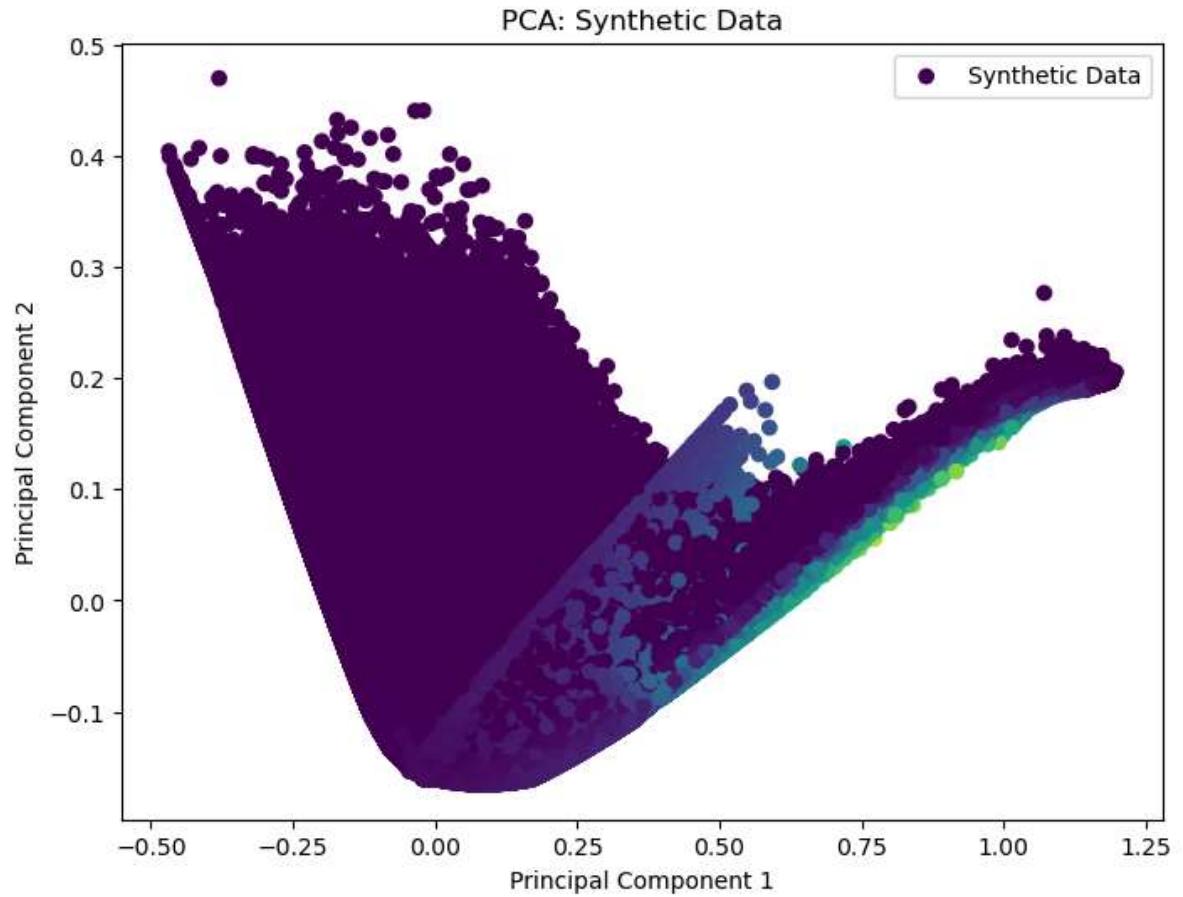
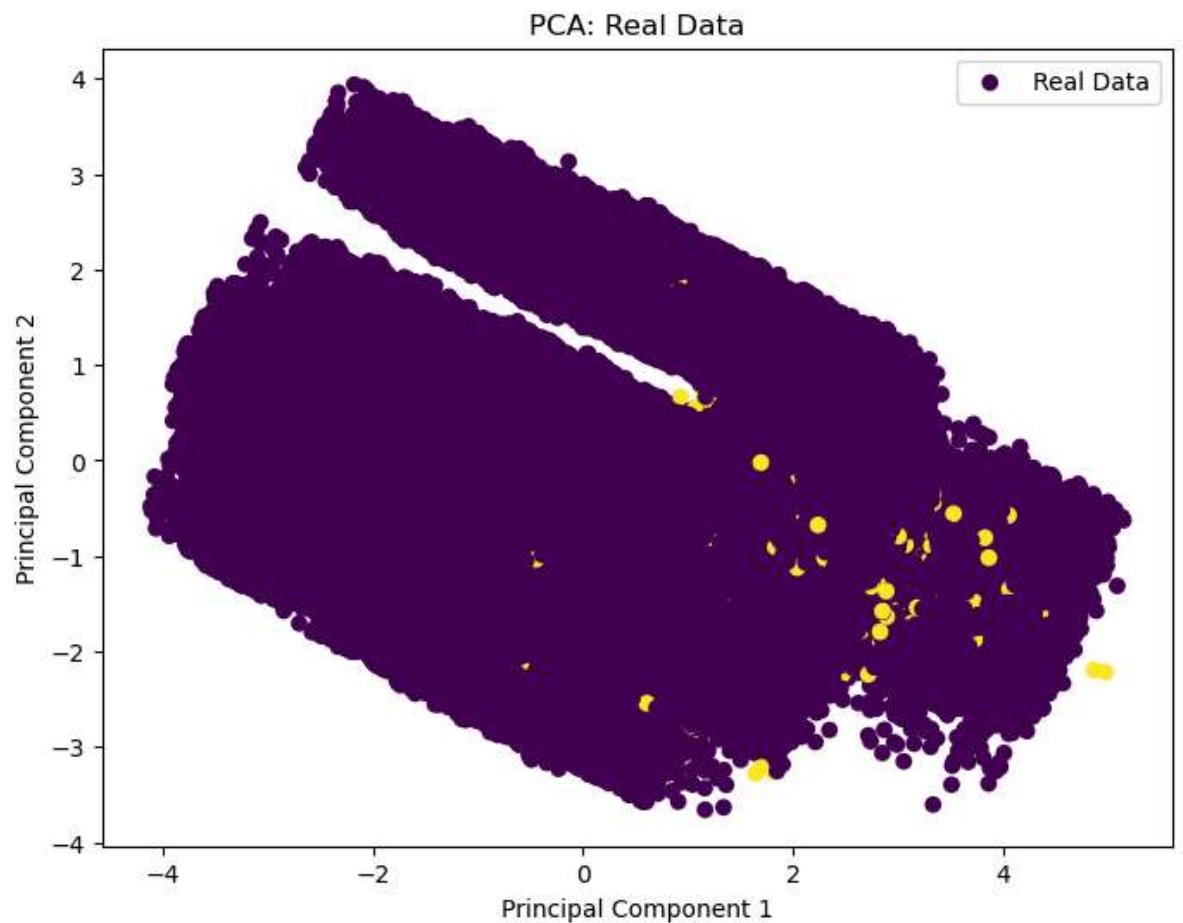
```
In [123]: 1 # Concatenate real and synthetic data
2 total_data = pd.concat([data, denormalizedDf_Labeled])
3
4 # Calculate the total number of duplicate rows in the concatenated data
5 dup_total = len(total_data) - len(total_data.drop_duplicates())
6
7 # Calculate the number of duplicate rows in the real data and synthetic data
8 dup_real = len(data) - len(data.drop_duplicates())
9 dup_synthetic = len(denormalizedDf_Labeled) - len(denormalizedDf_Labeled.drop_duplicates())
10
11 # Calculate the number of copies (duplicates not present in real or synthetic data)
12 copies = dup_total - dup_real - dup_synthetic
13
14 # Calculate the percentage of data that is a direct copy of the real data
15 copy_percentage = (copies / len(total_data)) * 100
16
17 # Calculate the percentage of self-copy (duplicate rows) in the synthetic data
18 self_copy_percentage = ((len(denormalizedDf_Labeled) - len(denormalizedDf_Labeled.drop_duplicates())))
19
20 print(f"Percentage of data that is a direct copy of the real data: {copy_percentage:.2f}%")
21 print(f"Percentage of data that is a self copy (duplicate rows): {self_copy_percentage:.2f}%")
```

Percentage of data that is a direct copy of the real data: 0.00%  
 Percentage of data that is a self copy (duplicate rows): 37.97%

## 2. Medium grained

- The "Medium grained" evaluation metric compares the distributions of real and synthetic data. This evaluation provides insights into how well the synthetic data replicates the distribution patterns of the real data across various features. By visualizing the distributions side by side, it helps assess the overall similarity or dissimilarity between the two datasets in terms of their statistical properties. This comparison can be crucial for understanding whether the synthetic data generation process accurately captures the underlying characteristics of the real data distribution.

```
In [136]:  
1 from sklearn.decomposition import PCA  
2 import matplotlib.pyplot as plt  
3  
4 X_PCA_real = X_train_scaled.iloc[:, :-1]  
5 y_PCA_real = X_train_scaled.iloc[:, -1]  
6 # Initialize PCA for real data  
7 pca_real = PCA(n_components=2)  
8 # Fit and transform the real data  
9 pca_real_result = pca_real.fit_transform(X_PCA_real)  
10  
11 X_PCA_syn = df_gen.iloc[:, :-1]  
12 y_PCA_syn = df_gen.iloc[:, -1]  
13  
14 # Initialize PCA for synthetic data  
15 pca_synthetic = PCA(n_components=2)  
16 # Fit and transform the synthetic data  
17 pca_synthetic_result = pca_synthetic.fit_transform(X_PCA_syn)  
18  
19 # Plot PCA results for real data  
20 plt.figure(figsize=(8, 6))  
21 plt.scatter(pca_real_result[:, 0], pca_real_result[:, 1], c=y_PCA_real, label=y_PCA_real)  
22 plt.xlabel('Principal Component 1')  
23 plt.ylabel('Principal Component 2')  
24 plt.title('PCA: Real Data')  
25 plt.legend()  
26 plt.show()  
27  
28 # Plot PCA results for synthetic data  
29 plt.figure(figsize=(8, 6))  
30 plt.scatter(pca_synthetic_result[:, 0], pca_synthetic_result[:, 1], c=y_PCA_syn, label=y_PCA_syn)  
31 plt.xlabel('Principal Component 1')  
32 plt.ylabel('Principal Component 2')  
33 plt.title('PCA: Synthetic Data')  
34 plt.legend()  
35 plt.show()  
36  
37
```



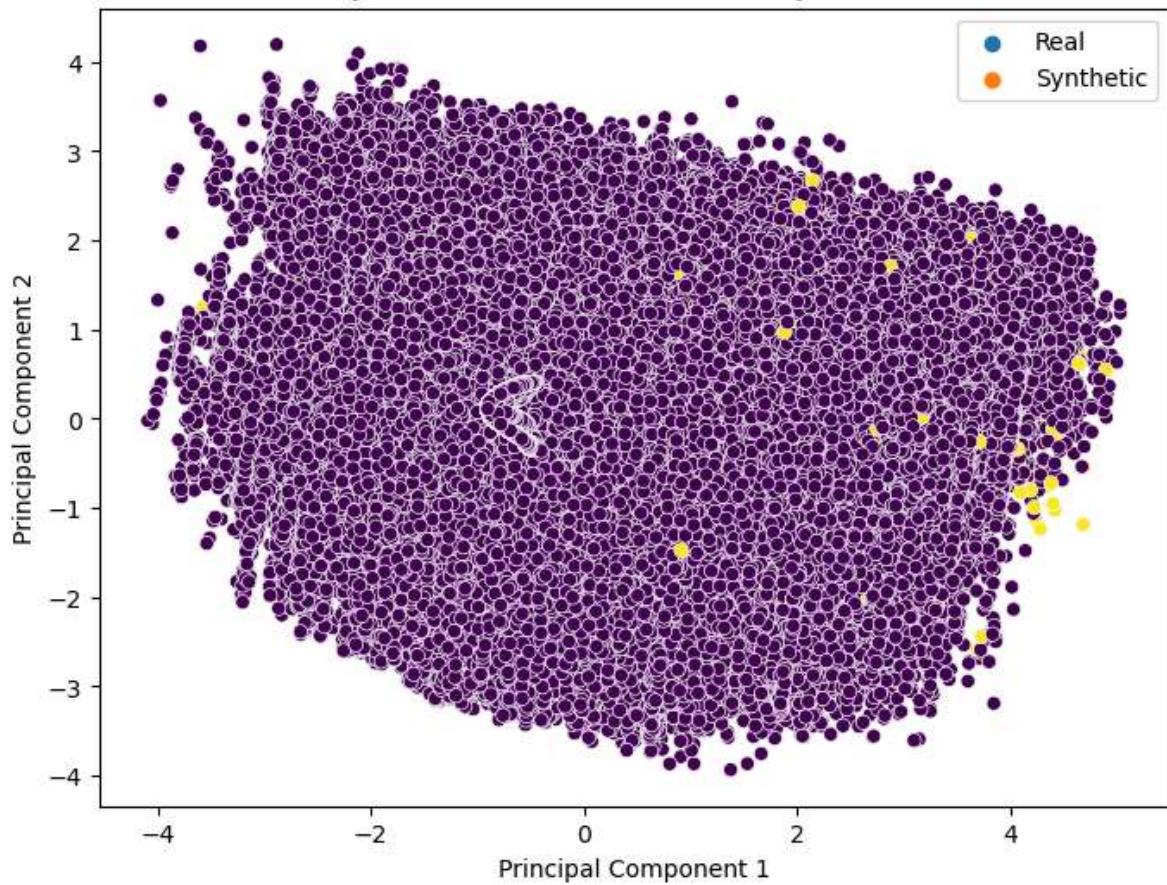
### 3. Fine grained

- The "Fine grained" evaluation metric involves comparing the joined distributions of real and synthetic data. This metric goes beyond comparing individual feature distributions and focuses on the overall joint distribution of multiple features. By analyzing how well the joint distributions align between real and synthetic data, this evaluation provides a deeper understanding of the data generation process's fidelity. It helps identify any discrepancies or patterns that may not be apparent when examining individual features separately, thus offering a more comprehensive assessment of the synthetic data's quality and resemblance to real-world data.

In [133]:

```
1 from sklearn.decomposition import PCA
2 import seaborn as sns
3 import matplotlib.pyplot as plt
4
5 # Combine real and synthetic data for comparison
6 combined_data = pd.concat([X_train_scaled, df_gen], axis=0)
7
8 combined_data_X = combined_data.iloc[:, :-1]
9 combined_data_y = combined_data.iloc[:, -1]
10
11 # Initialize PCA for combined data
12 pca_combined = PCA(n_components=2)
13 # Fit and transform the combined data
14 pca_combined_result = pca_combined.fit_transform(combined_data_X)
15
16 # Create a DataFrame for the PCA results
17 pca_df = pd.DataFrame(data=pca_combined_result, columns=['PC1', 'PC2'])
18 pca_df['Data Type'] = ['Real' for _ in range(len(X_train_scaled.iloc[:, :-1]))]
19
20 # Plot joint distribution of real and synthetic data
21 plt.figure(figsize=(8, 6))
22 sns.scatterplot(data=pca_df, c=combined_data_y, x='PC1', y='PC2', hue='Data Type')
23 plt.xlabel('Principal Component 1')
24 plt.ylabel('Principal Component 2')
25 plt.title('PCA: Joint Distribution of Real and Synthetic Data')
26 plt.legend()
27 plt.show()
```

## PCA: Joint Distribution of Real and Synthetic Data



In [ ]:

1