

# Data Preprocessing

```
In [1]: 1 import pandas as pd  
2 import numpy as np  
3 import os  
4 from sklearn.preprocessing import StandardScaler  
5 from sklearn.model_selection import train_test_split
```

```
In [2]: 1 def get_csv_files(target_dir):  
2     """Returns a list of all CSV file paths in the specified directory."""  
3     return [os.path.join(target_dir, f) for f in os.listdir(target_dir) if f.endswith('.csv')]  
4  
5 def save_plot(file_name, dir_path, plt):  
6     """Save the given plot to the specified directory."""  
7  
8     # Construct the full path for the image  
9     full_path = os.path.join(dir_path, file_name)  
10  
11    # Save the plot  
12    plt.savefig(full_path)  
13  
14    # Print a message indicating the image has been saved  
15    print(f"Image has been saved at {full_path}")
```

```
In [3]: 1 # Constants  
2 DTYPES = {"cashier": "object", "routing_no": "object", "tm_disc": "object"}  
3 DROP_COLS = ["rcpt_id", 'tndr_acct', 'routing_no', 'start_tm', 'end_tm', 'sls_d', 'crt_whl_a']  
4 FLAG_COLS = ['gc_flg', 'prc_ovrd_exc_flg', 'void_flag', 'tm_flg']  
5 FLAG_REPLACE = {'Y': 1, '': 0, ' ': 0}
```

```
In [4]: 1 Image_Folder = 'Images'  
2 Model_Folder = 'models'  
3 Result_Folder = 'Results'  
4 Data_Folder = 'data'  
5 Data_SubFolder = 'Target_Data\Receipts_2'  
6 SubFolder = "All_Purchases"  
7  
8 data_dir_path = os.path.join(Data_Folder, Data_SubFolder)  
9  
10 image_dir_path = os.path.join(Image_Folder, SubFolder)  
11 os.makedirs(image_dir_path, exist_ok=True)  
12  
13 model_dir_path = os.path.join(Model_Folder, SubFolder)  
14 os.makedirs(model_dir_path, exist_ok=True)  
15  
16 result_dir_path = os.path.join(Result_Folder, SubFolder)  
17 os.makedirs(result_dir_path, exist_ok=True)
```

```
In [5]: 1 # Get CSV files  
2 csv_files = get_csv_files(data_dir_path)  
3  
4 # Import and concatenate CSV files  
5 dfs = [pd.read_csv(file, dtype=DTYPES) for file in csv_files]  
6 df = pd.concat(dfs, ignore_index=True)  
7  
8 # Backup original dataframe  
9 df_orig = df.copy().reset_index(drop=True)
```

```
In [6]: 1 # Manipulate the 'register_i' column
2 df['register_i'] = df['rcpt_id'].str.split('-').str[2].str.lstrip('0')
3 df['register_i'] = df['register_i'].replace('', 'Unknown').fillna('Unknown')
4
5 # Convert columns to binary representations
6 df['guest_profile_id'] = df['guest_profile_id'].notna().astype(int)
7 df['circle_id'] = df['circle_id'].notna().astype(int)
8
9 # Update FLAG_COLS and other columns
10 df = (df.fillna({col: 0 for col in FLAG_COLS})
11     .replace({col: FLAG_REPLACE for col in FLAG_COLS})
12     .astype({col: int for col in FLAG_COLS})
13     .assign(tndr_type=lambda x: x['tndr_type'].str.strip(),
14             sls_d=lambda x: pd.to_datetime(x['sls_d'], format='%Y-%m-%d'),
15             start_tm=lambda x: pd.to_datetime(x['sls_d'].dt.date.astype(str) + ' ' + x['start_tm'].astype(str)))
16     .reset_index(drop=True))
17
18 # Extract date and time features
19 df['year'] = df['start_tm'].dt.year
20 df['month'] = df['start_tm'].dt.month
21 df['hour'] = df['start_tm'].dt.hour
22 df['day_of_year'] = df['start_tm'].dt.dayofyear
23 df['day_of_month'] = df['start_tm'].dt.day
24 df['day_of_week'] = df['start_tm'].dt.dayofweek
25
26 # Further preprocess
27 df['tndr_type'] = df['tndr_type'].replace('', 'Unknown').fillna('Unknown')
28 df['cshbk_amt'] = df['cshbk_amt'].fillna(0)
29 df['tm_disc'] = df['tm_disc'].fillna(0)
30
31 # Fill NaN values
32 df = df.fillna(0)
33 for column in df.select_dtypes(include=['object']).columns:
34     if column not in FLAG_COLS:
35         df[column] = df[column].replace('0', 'Unknown')
36
37 df.loc[df['register_i'].isin([61, 62, 63, 64]), 'cashier'] = 'Self_Checkout'
```

```
In [7]: 1 from sklearn.model_selection import train_test_split
2 from sklearn.preprocessing import MinMaxScaler
3
4 #Drop unnecessary columns
5 df.drop(columns=DROP_COLS, inplace=True)
6
7 # Get columns that are of 'object' dtype
8 columns_to_dummy = df.select_dtypes(include=['object']).columns
9
10 # Create dummy variables for the selected columns
11 df = pd.get_dummies(df, columns=columns_to_dummy).astype(int)
12
13 # Split data into training and validation sets (80% train, 20% validation)
14 train_data, val_data = train_test_split(df, test_size=0.2)
15
16 scaler = MinMaxScaler(feature_range=(-1, 1))
17 train_data_scaled = scaler.fit_transform(train_data)
18 val_data_scaled = scaler.transform(val_data)
19
20 train_data_scaled.shape, val_data_scaled.shape
```

Out[7]: ((928936, 735), (232234, 735))

## Neural Networks

```
In [8]: 1 import torch
2 import torch.nn as nn
3 import torch.optim as optim
4 import torch.nn.functional as F
5 from torch.utils.data import DataLoader, TensorDataset
```

```
In [9]: 1 # Hyperparameters
2 EPOCHS = 20
3 BATCH_SIZE = 128
4 LEARNING_RATE = 0.001
5 LOAD_MODEL = True
6 SAVE_MODEL = True
7
8 models = {}
9
10 device = torch.device("cuda" if torch.cuda.is_available() else "cpu")
```

```
In [10]: 1 # Convert data to PyTorch tensors
2 train_dataset = TensorDataset(torch.tensor(train_data_scaled, dtype=torch.float32))
3 train_dataloader = DataLoader(train_dataset, batch_size=BATCH_SIZE, shuffle=True)
4 val_dataset = TensorDataset(torch.tensor(val_data_scaled, dtype=torch.float32))
5 val_dataloader = DataLoader(val_dataset, batch_size=BATCH_SIZE, shuffle=False)
```

## Train/Val Functions

```
In [11]: 1 from timeit import default_timer as timer
2 from tqdm.auto import tqdm
3 import matplotlib.pyplot as plt
```

```
In [12]: 1 def train_step(model: torch.nn.Module,
2                  dataloader: torch.utils.data.DataLoader,
3                  loss_fn: torch.nn.Module,
4                  optimizer: torch.optim.Optimizer,
5                  device=device,
6                  sparse: bool = False):
7
8     # Switch the model to train mode
9     model.train()
10
11    # Initialize variables to store total training Loss
12    train_loss = 0
13
14    # Loop over batches from the DataLoader
15    for batch, X in enumerate(dataloader):
16
17        # Send batch of images to the computation device (CPU/GPU)
18        X = X[0].to(device)
19
20        # Perform a forward pass through the model to get the predictions
21        outputs = model(X)
22
23        # Calculate the loss between the predictions and actual values
24        if sparse:
25            X_pred = outputs
26            loss = loss_fn(model, X_pred, X)
27        else:
28            X_pred = outputs
29            loss = loss_fn(X_pred, X)
30
31        # Add up the loss values
32        train_loss += loss.item()
33
34        # Reset the gradients from the previous iteration
35        optimizer.zero_grad()
36
37        # Perform backward propagation to calculate gradients
38        loss.backward()
39
40        # Perform a step of optimization
41        optimizer.step()
42
43    train_loss /= len(dataloader) # Calculate average training Loss
44    return train_loss # Return average training Loss
```

In [13]:

```
1 # Create eval_step() function
2 def val_step(model: torch.nn.Module,
3             dataloader: torch.utils.data.DataLoader,
4             loss_fn: torch.nn.Module,
5             device=device,
6             sparse: bool = False):
7
8     # Switches the model to evaluation mode
9     model.eval()
10
11    # Initialize variables to store total test loss and accuracy
12    val_loss = 0
13
14    # Disable calculation of gradients for performance boost during inference
15    with torch.inference_mode():
16
17        # Loop over batches from the DataLoader
18        for batch, X in enumerate(dataloader):
19
20            # Send batch of images and labels to the computation device (CPU/GPU)
21            X = X[0].to(device)
22
23            # Perform a forward pass through the model to get the predictions
24            outputs = model(X)
25
26            # Calculate the Loss between the predictions and actual values
27            if sparse:
28                X_pred = outputs
29                loss = loss_fn(model, X_pred, X)
30            else:
31                X_pred = outputs
32                loss = loss_fn(X_pred, X)
33            # Add up the loss values
34            val_loss += loss.item()
35
36    val_loss /= len(dataloader) # Calculate average val loss
37    return val_loss # Return average val loss
```

In [14]:

```

1 # Create train() function
2 def train(model_name: str,
3           model: torch.nn.Module,
4           train_dataloader: torch.utils.data.DataLoader,
5           val_dataloader: torch.utils.data.DataLoader,
6           optimizer: torch.optim.Optimizer,
7           loss_fn: torch.nn.Module,
8           epochs: int = 5,
9           device = device,
10          sparse: bool = False,
11          load_model: bool = True,
12          save_model: bool = True,
13          model_dir_path = model_dir_path):
14
15     # Define the full model path including the file name
16     model_path = os.path.join(model_dir_path, f"{model_name}.pth")
17
18     # Initialize a dictionary to store training and validation losses and accuracies for each epoch
19     results = {"train_loss": [], "val_loss": []}
20
21     # Load the existing model and losses if required
22     if load_model:
23         loaded = load_model_fn(model, model_path, results)
24         if loaded:
25             return results
26
27
28     # Establish the start time for training
29     start_time = timer()
30
31     # Loop over epochs
32     for epoch in tqdm(range(epochs)):
33         # Execute a training step and get training loss
34         train_loss = train_step(model=model,
35                                dataloader=train_dataloader,
36                                loss_fn=loss_fn,
37                                optimizer=optimizer,
38                                device=device,
39                                sparse=sparse)
34         # Execute a validation step and get validation loss
35         val_loss = val_step(model=model,
36                             dataloader=val_dataloader,
37                             loss_fn=loss_fn,
38                             device=device,
39                             sparse=sparse)
40
41         # Print Losses and accuracies for this epoch
42         print(f"Epoch: {epoch} | Train loss: {train_loss:.4f} | Val loss: {val_loss:.4f}")
43
44         # Append losses and accuracies to results dictionary
45         results["train_loss"].append(train_loss)
46         results["val_loss"].append(val_loss)
47
48     # Establish the end time for training
49     end_time = timer()
50     train_time = end_time - start_time
51     results['train_time'] = train_time
52
53     if save_model:
54         save_model_fn(model, model_path, results)
55
56     return results

```

In [15]:

```

1 def load_model_fn(model, model_path, results):
2     if os.path.isfile(model_path):
3         checkpoint = torch.load(model_path)
4         model.load_state_dict(checkpoint['state_dict'])
5         results['train_loss'] = checkpoint['train_loss']
6         results['val_loss'] = checkpoint['val_loss']
7         results['train_time'] = checkpoint['train_time']
8         print(f"Model loaded from {model_path}")
9         return True
10    else:
11        print(f"No existing model found at {model_path}. Training a new model.")
12        return False

```

```
In [16]: 1 def save_model_fn(model, model_path, results):
2     save_dict = {
3         'state_dict': model.state_dict(),
4         'train_loss': results['train_loss'],
5         'val_loss': results['val_loss'],
6         'train_time': results['train_time']
7     }
8     torch.save(save_dict, model_path)
9     print(f"Model saved at {model_path}")
```

```
In [17]: 1 def plot_train_val_loss(results, model, image_dir_path=image_dir_path):
2     plt.figure(figsize=(10,5))
3
4     plt.plot(results['train_loss'], label='Training Loss')
5     plt.plot(results['val_loss'], label='Validation Loss')
6
7     plt.xlabel('Epoch')
8     plt.ylabel('Loss')
9     plt.title(f'Test and Validation Loss for {model}')
10    plt.legend()
11
12    save_plot(f'{model}_Train_Val_Loss", image_dir_path, plt)
13
14    plt.show()
```

## Model 1

Regular Autoencoder

```
In [18]: 1 # Define the autoencoder model
2 class Autoencoder(torch.nn.Module):
3     def __init__(self, input_dim):
4         super(Autoencoder, self).__init__()
5         self.encoder = nn.Sequential(
6             nn.Linear(input_dim, 256),
7             nn.ReLU(),
8
9             nn.Linear(256, 128),
10            nn.ReLU(),
11
12            nn.Linear(128, 64),
13            nn.ReLU(),
14
15            nn.Linear(64, 32),
16            nn.ReLU(),
17
18            nn.Linear(32, 16),
19            nn.ReLU(),
20        )
21        self.decoder = nn.Sequential(
22            nn.Linear(16, 32),
23            nn.ReLU(),
24
25            nn.Linear(32, 64),
26            nn.ReLU(),
27
28            nn.Linear(64, 128),
29            nn.ReLU(),
30
31            nn.Linear(128, 256),
32            nn.ReLU(),
33
34            nn.Linear(256, input_dim),
35            nn.Tanh(),
36        )
37
38    def forward(self, x):
39        x = self.encoder(x)
40        x = self.decoder(x)
41        return x
```

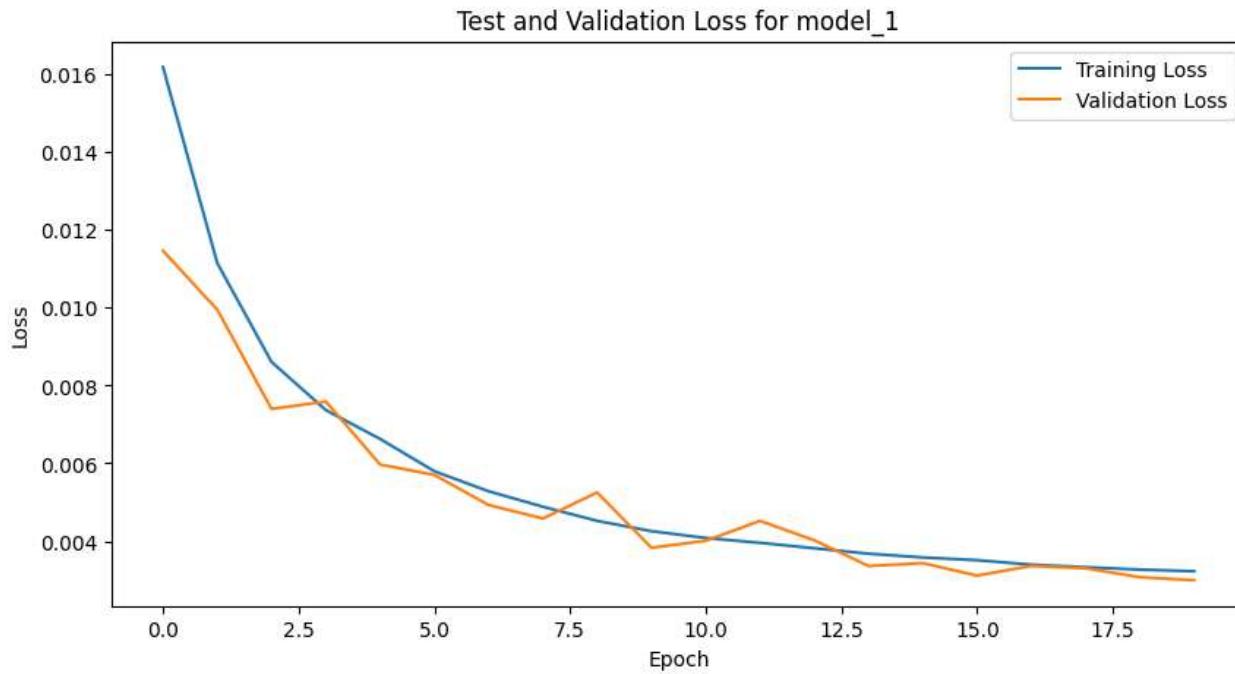
```
In [19]: 1 # Initialize the model, loss function and optimizer
2 model_1 = Autoencoder(input_dim=train_data_scaled.shape[1])
3 loss_fn_1 = nn.MSELoss()
4 optimizer_1 = optim.Adam(model_1.parameters(), lr=LEARNING_RATE)
5
6 # Move model to appropriate device
7 model_1.to(device)
8 models["model_1"] = model_1
```

```
In [20]: 1 model_1_results = train(model_name = list(models.keys())[0],
2                               model=model_1,
3                               train_dataloader=train_dataloader,
4                               val_dataloader=val_dataloader,
5                               optimizer=optimizer_1,
6                               loss_fn=loss_fn_1,
7                               epochs=EPOCHS,
8                               device=device,
9                               load_model=LOAD_MODEL,
10                              save_model=SAVE_MODEL)
```

Model loaded from models\All\_Purchases\model\_1.pth

```
In [21]: 1 plot_train_val_loss(model_1_results, list(models.keys())[0])
```

Image has been saved at Images\All\_Purchases\model\_1\_Train\_Val\_Loss



## Model 2

Deeper autoencoder with dropout and normalization

```
In [22]: 1 # Define the autoencoder model
2 import torch.nn as nn
3
4 class DeepAutoencoder(torch.nn.Module):
5     def __init__(self, input_dim):
6         super(DeepAutoencoder, self).__init__()
7
8         # Encoder
9         self.encoder = nn.Sequential(
10             nn.Linear(input_dim, 512),
11             nn.ReLU(),
12             nn.BatchNorm1d(512),
13             nn.Dropout(0.5),
14
15             nn.Linear(512, 256),
16             nn.ReLU(),
17             nn.BatchNorm1d(256),
18             nn.Dropout(0.5),
19
20             nn.Linear(256, 128),
21             nn.ReLU(),
22             nn.BatchNorm1d(128),
23             nn.Dropout(0.5),
24
25             nn.Linear(128, 64),
26             nn.ReLU(),
27             nn.BatchNorm1d(64),
28
29             nn.Linear(64, 32),
30             nn.ReLU(),
31             nn.BatchNorm1d(32),
32
33             nn.Linear(32, 16),
34             nn.ReLU(),
35             nn.BatchNorm1d(16)
36         )
37
38         # Decoder
39         self.decoder = nn.Sequential(
40             nn.Linear(16, 32),
41             nn.ReLU(),
42             nn.BatchNorm1d(32),
43
44             nn.Linear(32, 64),
45             nn.ReLU(),
46             nn.BatchNorm1d(64),
47
48             nn.Linear(64, 128),
49             nn.ReLU(),
50             nn.BatchNorm1d(128),
51
52             nn.Linear(128, 256),
53             nn.ReLU(),
54             nn.BatchNorm1d(256),
55
56             nn.Linear(256, input_dim),
57             nn.Tanh()
58         )
59
60     def forward(self, x):
61         x = self.encoder(x)
62         x = self.decoder(x)
63
64         return x
```

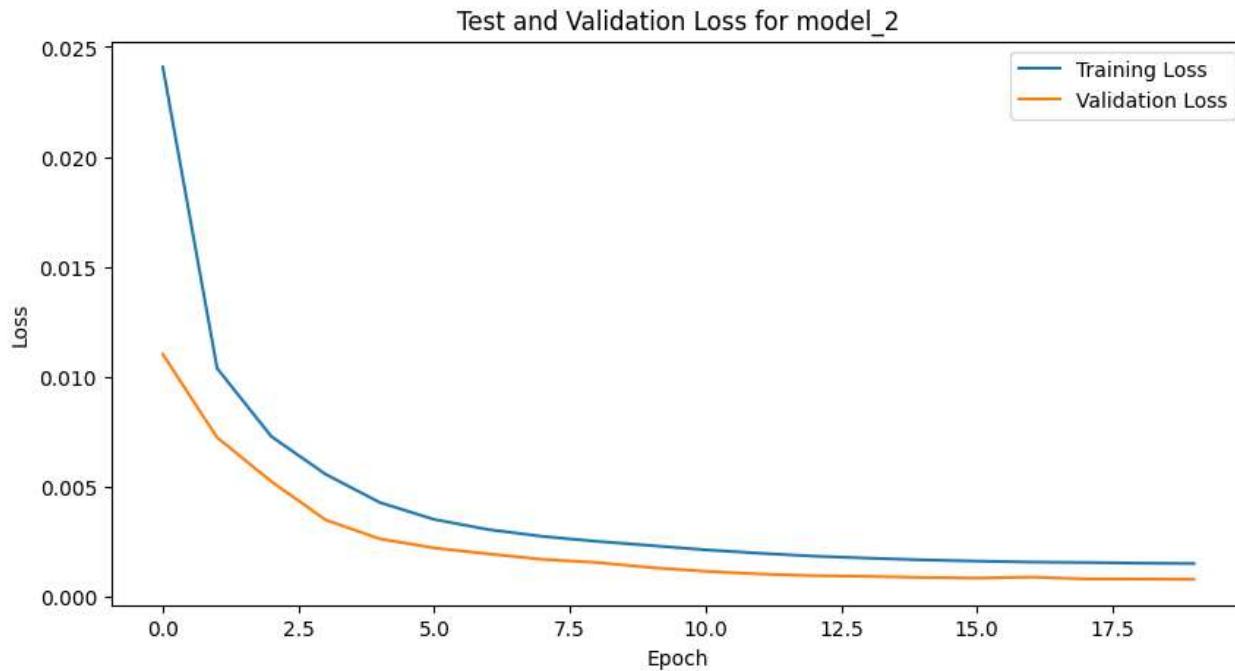
```
In [23]: 1 # Initialize the model, loss function and optimizer
2 model_2 = DeepAutoencoder(input_dim=train_data_scaled.shape[1])
3 loss_fn_2 = nn.MSELoss()
4 optimizer_2 = optim.Adam(model_2.parameters(), lr=LEARNING_RATE)
5
6 # Move model to appropriate device
7 model_2.to(device)
8 models["model_2"] = model_2
```

```
In [24]: 1 model_2_results = train(model_name = list(models.keys())[1],
2                               model=model_2,
3                               train_dataloader=train_dataloader,
4                               val_dataloader=val_dataloader,
5                               optimizer=optimizer_2,
6                               loss_fn=loss_fn_2,
7                               epochs=EPOCHS,
8                               device=device,
9                               load_model=LOAD_MODEL,
10                              save_model=SAVE_MODEL)
```

Model loaded from models\All\_Purchases\model\_2.pth

```
In [25]: 1 plot_train_val_loss(model_2_results, list(models.keys())[1])
```

Image has been saved at Images\All\_Purchases\model\_2\_Train\_Val\_Loss



## Model 3

Denoising Autoencoder

In [26]:

```

1 class DenoisingAutoencoder(nn.Module):
2     def __init__(self, input_dim, noise_factor=0.5):
3         super(DenoisingAutoencoder, self).__init__()
4         self.noise_factor = noise_factor
5         # Encoder
6         self.encoder = nn.Sequential(
7             nn.Linear(input_dim, 512),
8             nn.ReLU(),
9             nn.BatchNorm1d(512),
10            nn.Dropout(0.5),
11
12            nn.Linear(512, 256),
13            nn.ReLU(),
14            nn.BatchNorm1d(256),
15            nn.Dropout(0.5),
16
17            nn.Linear(256, 128),
18            nn.ReLU(),
19            nn.BatchNorm1d(128),
20            nn.Dropout(0.5),
21
22            nn.Linear(128, 64),
23            nn.ReLU(),
24            nn.BatchNorm1d(64),
25
26            nn.Linear(64, 32),
27            nn.ReLU(),
28            nn.BatchNorm1d(32),
29
30            nn.Linear(32, 16),
31            nn.ReLU(),
32            nn.BatchNorm1d(16)
33        )
34
35         # Decoder
36         self.decoder = nn.Sequential(
37             nn.Linear(16, 32),
38             nn.ReLU(),
39             nn.BatchNorm1d(32),
40
41             nn.Linear(32, 64),
42             nn.ReLU(),
43             nn.BatchNorm1d(64),
44
45             nn.Linear(64, 128),
46             nn.ReLU(),
47             nn.BatchNorm1d(128),
48
49             nn.Linear(128, 256),
50             nn.ReLU(),
51             nn.BatchNorm1d(256),
52
53             nn.Linear(256, input_dim),
54             nn.Tanh()
55         )
56
57     def forward(self, x):
58         x_noisy = x + self.noise_factor * torch.randn_like(x) # Add noise to the input
59         x_noisy = torch.clamp(x_noisy, 0., 1.) # Ensure the noisy input values are between 0 and 1.
60         x_encoded = self.encoder(x_noisy) # Encode the noisy input
61         x_decoded = self.decoder(x_encoded) # Decode the encoded representation
62         return x_decoded

```

In [27]:

```

1 # Initialize the model, loss function and optimizer
2 model_3 = DenoisingAutoencoder(input_dim=train_data_scaled.shape[1])
3 loss_fn_3 = nn.MSELoss()
4 optimizer_3 = optim.Adam(model_3.parameters(), lr=LEARNING_RATE)
5
6 # Move model to appropriate device
7 model_3.to(device)
8 models["model_3"] = model_3

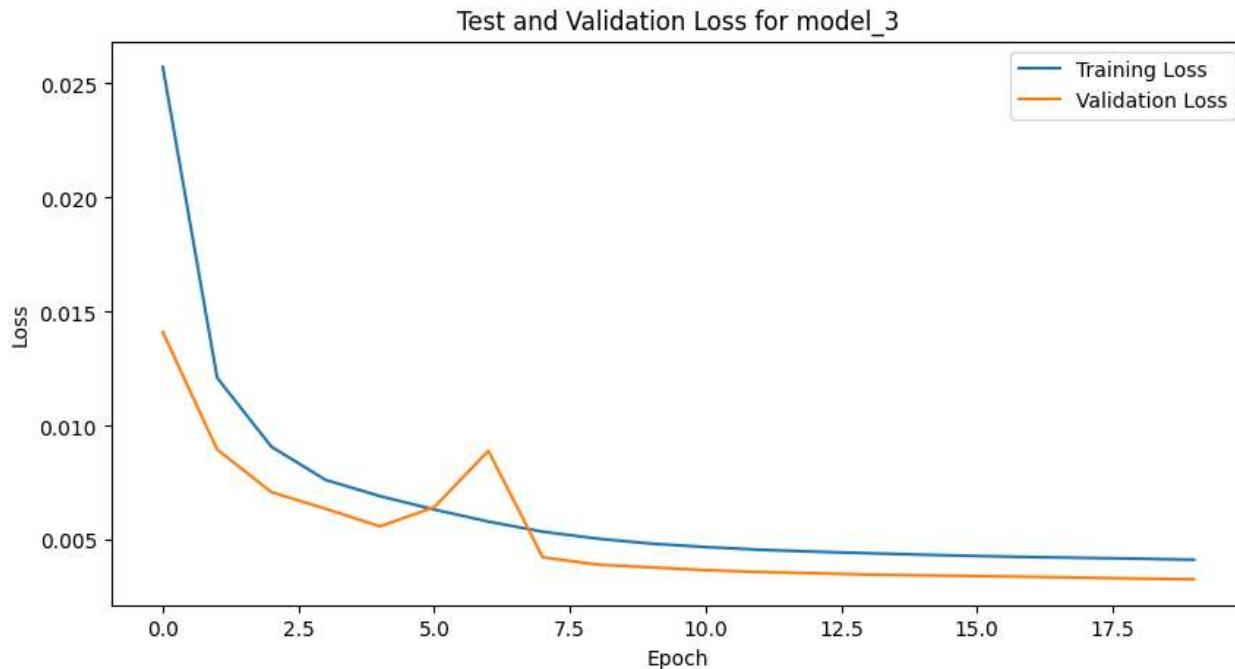
```

```
In [28]: 1 model_3_results = train(model_name = list(models.keys())[2],
2                               model=model_3,
3                               train_dataloader=train_dataloader,
4                               val_dataloader=val_dataloader,
5                               optimizer=optimizer_3,
6                               loss_fn=loss_fn_3,
7                               epochs=EPOCHS,
8                               device=device,
9                               load_model=LOAD_MODEL,
10                              save_model=SAVE_MODEL)
```

Model loaded from models\All\_Purchases\model\_3.pth

```
In [29]: 1 plot_train_val_loss(model_3_results, list(models.keys())[2])
```

Image has been saved at Images\All\_Purchases\model\_3\_Train\_Val\_Loss



## Model 4

Sparse Autoencoder

In [30]:

```

1 class SparseAutoencoder(nn.Module):
2     def __init__(self, input_dim):
3         super(SparseAutoencoder, self).__init__()
4         # Encoder
5         self.encoder = nn.Sequential(
6             nn.Linear(input_dim, 512),
7             nn.ReLU(),
8             nn.BatchNorm1d(512),
9             nn.Dropout(0.5),
10
11            nn.Linear(512, 256),
12            nn.ReLU(),
13            nn.BatchNorm1d(256),
14            nn.Dropout(0.5),
15
16            nn.Linear(256, 128),
17            nn.ReLU(),
18            nn.BatchNorm1d(128),
19            nn.Dropout(0.5),
20
21            nn.Linear(128, 64),
22            nn.ReLU(),
23            nn.BatchNorm1d(64),
24
25            nn.Linear(64, 32),
26            nn.ReLU(),
27            nn.BatchNorm1d(32),
28
29            nn.Linear(32, 16),
30            nn.ReLU(),
31            nn.BatchNorm1d(16)
32        )
33
34         # Decoder
35         self.decoder = nn.Sequential(
36             nn.Linear(16, 32),
37             nn.ReLU(),
38             nn.BatchNorm1d(32),
39
40             nn.Linear(32, 64),
41             nn.ReLU(),
42             nn.BatchNorm1d(64),
43
44             nn.Linear(64, 128),
45             nn.ReLU(),
46             nn.BatchNorm1d(128),
47
48             nn.Linear(128, 256),
49             nn.ReLU(),
50             nn.BatchNorm1d(256),
51
52             nn.Linear(256, input_dim),
53             nn.Tanh()
54        )
55
56     def forward(self, x):
57         x = self.encoder(x)
58         self.mean_activation = torch.mean(x, dim=0) # Compute the mean activation
59         x = self.decoder(x)
60         return x

```

In [31]:

```

1 def sparse_loss_fn(model, X, X_pred, beta=1e-3, rho=0.05, epsilon=1e-7):
2     # Reconstruction loss
3     mse_loss = nn.MSELoss()(X_pred, X)
4
5     # Ensure mean_activation is not exactly 0 or 1
6     mean_activation = torch.clamp(model.mean_activation, min=epsilon, max=1-epsilon)
7
8     # Sparsity loss
9     sparsity_loss = torch.abs(torch.sum(rho * torch.log(rho / mean_activation) + (1 - rho) * torch.log((1 - rho) / (1 - mean_activation))))
10
11    # Total loss
12    total_loss = mse_loss + beta * sparsity_loss
13
14    return total_loss

```

```
In [32]: 1 # Initialize the model, loss function and optimizer
2 model_4 = SparseAutoencoder(input_dim=train_data_scaled.shape[1])
3 loss_fn_4 = sparse_loss_fn
4 optimizer_4 = optim.Adam(model_4.parameters(), lr=LEARNING_RATE)
5
6 # Move model to appropriate device
7 model_4.to(device)
8 models["model_4"] = model_4
```

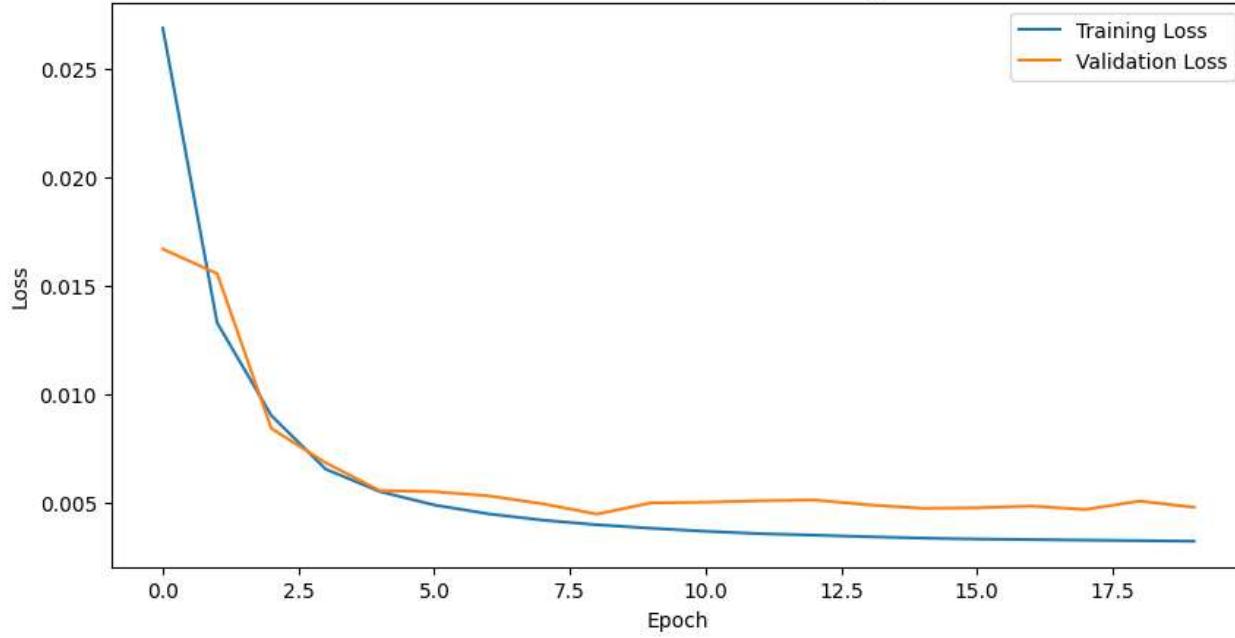
```
In [33]: 1 model_4_results = train(model_name = list(models.keys())[3],
2                               model=model_4,
3                               train_dataloader=train_dataloader,
4                               val_dataloader=val_dataloader,
5                               optimizer=optimizer_4,
6                               loss_fn=loss_fn_4,
7                               epochs=EPOCHS,
8                               device=device,
9                               sparse=True,
10                              load_model=LOAD_MODEL,
11                              save_model=SAVE_MODEL)
```

Model loaded from models\All\_Purchases\model\_4.pth

```
In [34]: 1 plot_train_val_loss(model_4_results, list(models.keys())[3])
```

Image has been saved at Images\All\_Purchases\model\_4\_Train\_Val\_Loss

Test and Validation Loss for model\_4



## Comparison Graphs

In [35]:

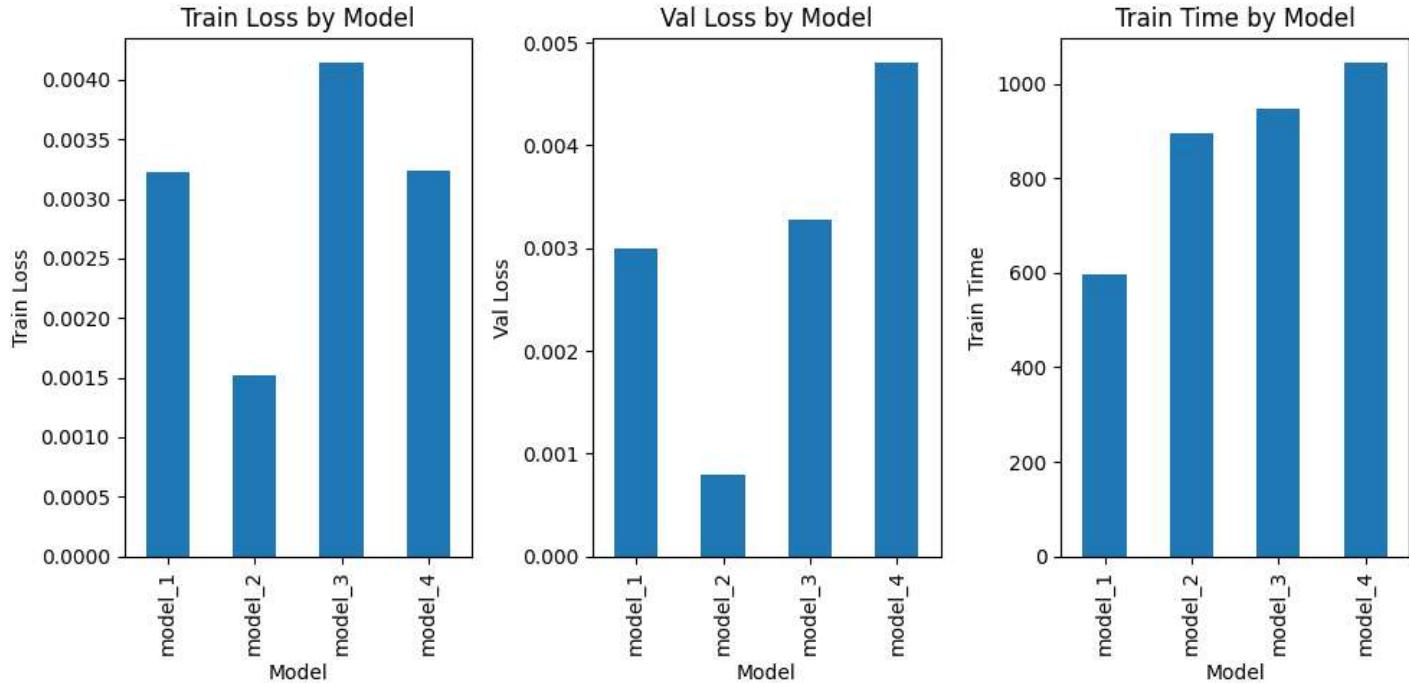
```
1 # Create a dictionary of model results
2 model_dict = {
3     'model_1': model_1_results,
4     'model_2': model_2_results,
5     'model_3': model_3_results,
6     'model_4': model_4_results
7 }
8
9 final_results = {
10     model_name: {
11         key: value[-1] if isinstance(value, list) and len(value) > 1 else value
12         for key, value in model.items()
13     }
14     for model_name, model in model_dict.items()
15 }
16
17 # Create DataFrame with correct orientation and add training time
18 results = pd.DataFrame(final_results).T
19
20 results
```

Out[35]:

	train_loss	val_loss	train_time
model_1	0.003226	0.002995	597.201953
model_2	0.001518	0.000798	896.574678
model_3	0.004138	0.003279	946.267649
model_4	0.003234	0.004800	1043.932900

```
In [36]: 1 plt.figure(figsize=(10,5))
2
3 plt.subplot(1, 3, 1)
4 results['train_loss'].plot(kind='bar')
5 plt.xlabel('Model')
6 plt.ylabel('Train Loss')
7 plt.title('Train Loss by Model')
8
9 plt.subplot(1, 3, 2)
10 results['val_loss'].plot(kind='bar')
11 plt.xlabel('Model')
12 plt.ylabel('Val Loss')
13 plt.title('Val Loss by Model')
14
15 plt.subplot(1, 3, 3)
16 results['train_time'].plot(kind='bar')
17 plt.xlabel('Model')
18 plt.ylabel('Train Time')
19 plt.title('Train Time by Model')
20
21 plt.tight_layout() # Adjusts spacing between subplots for better layout
22 file_name = "TrainLoss_ValLoss_TrainTime"
23
24 save_plot("TrainLoss_ValLoss_TrainTime_Bar", image_dir_path, plt)
25
26 plt.show()
```

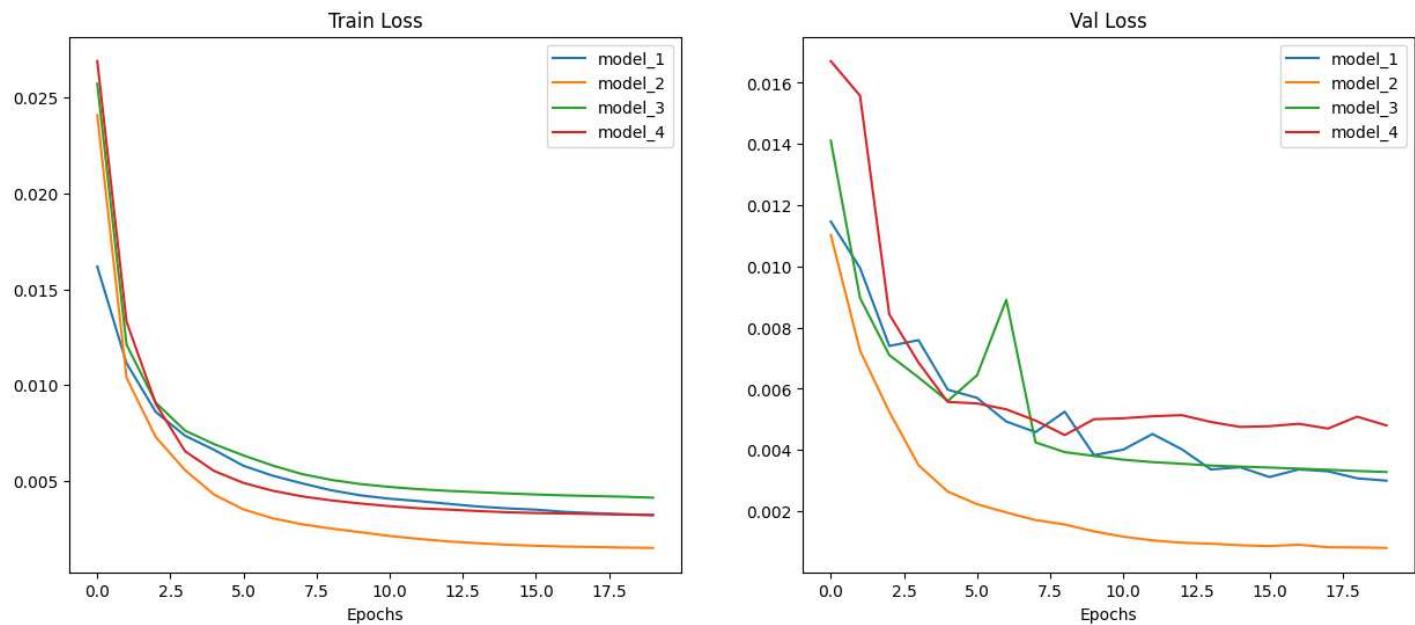
Image has been saved at Images\All\_Purchases\TrainLoss\_ValLoss\_TrainTime\_Bar



In [37]:

```
1 # Setup a plot
2 plt.figure(figsize=(15,6))
3
4 # Get number of epochs
5 epochs = range(len(model_1_results["train_loss"]))
6
7 # Plot train Loss
8 plt.subplot(1, 2, 1)
9 for i in range(len(model_dict)):
10     plt.plot(epochs, model_dict[list(model_dict)[i]]['train_loss'], label=list(model_dict)[i])
11 plt.title("Train Loss")
12 plt.xlabel("Epochs")
13 plt.legend()
14
15 # Plot val Loss
16 plt.subplot(1, 2, 2)
17 for i in range(len(model_dict)):
18     plt.plot(epochs, model_dict[list(model_dict)[i]]['val_loss'], label=list(model_dict)[i])
19 plt.title("Val Loss")
20 plt.xlabel("Epochs")
21 plt.legend()
22
23 save_plot("TrainLoss_ValLoss_Line", image_dir_path, plt)
24
25 plt.show()
```

Image has been saved at Images\All\_Purchases\TrainLoss\_ValLoss\_Line



## Graphing in Latent Space



In [38]:

```
1 import numpy as np
2 import matplotlib.pyplot as plt
3 from sklearn.decomposition import PCA
4 from sklearn.cluster import KMeans
5 from tqdm import tqdm
6 import torch
7 from torch import linalg as LA
8
9 def apply_pca_optimized(data, n_components=None, visualize=True, labels=None, image_dir_path=image_dir_path):
10     """
11         Apply Principal Component Analysis (PCA) on the given data and optionally visualize the results.
12         The data will be reduced to 2 dimensions for visualization.
13
14     :param data: Input data to apply PCA on.
15     :param n_components: Number of principal components to consider for transformation. If None, use all components.
16     :param visualize: Whether to visualize the results. Defaults to True.
17     :return: Transformed data after applying PCA.
18     """
19
20     # Apply PCA using sklearn
21     pca = PCA(n_components=n_components)
22     transformed_data = pca.fit_transform(data)
23
24     # Visualize the results if required
25     if visualize:
26         visualize_pca(transformed_data, labels, image_dir_path)
27
28     # Print the total explained variance
29     print("Total Explained Variance ({} components): {:.2%}".format(n_components if n_components else data.shape[1], np.sum
30
31     return transformed_data
32
33 def visualize_pca(transformed_data, labels=None, image_dir_path=image_dir_path):
34     """
35         Visualize the PCA results in 2D space.
36     # Apply PCA for 2D visualization
37     pca_2d = PCA(n_components=2)
38     transformed_data_2d = pca_2d.fit_transform(transformed_data)
39
40     plt.figure(figsize=(10,10))
41     plt.scatter(transformed_data_2d[:, 0], transformed_data_2d[:, 1], c=labels, s=1, cmap='viridis')
42     plt.xlabel("Principal Component 1")
43     plt.ylabel("Principal Component 2")
44     plt.title("2D PCA Projection")
45
46     if labels is not None:
47         save_plot("PCA_Visualization_Labels", image_dir_path, plt)
48     else:
49         save_plot("PCA_Visualization", image_dir_path, plt)
50
51     plt.show()
52
53 def plot_cumulative_variance(data, max_components=15, image_dir_path=image_dir_path):
54     """
55         Plot cumulative explained variance against the number of components.
56     # Convert data to PyTorch tensor and move to the device
57     data_tensor = torch.tensor(data, dtype=torch.float32).to(device)
58
59     # Center the data
60     mean_centered_data = data_tensor - data_tensor.mean(dim=0)
61
62     # Compute the covariance matrix
63     covariance_matrix = torch.mm(mean_centered_data.T, mean_centered_data) / (data_tensor.size(0) - 1)
64
65     # Compute eigenvalues
66     eigenvalues, _ = LA.eig(covariance_matrix)
67     eigenvalues = torch.real(eigenvalues)
68
69     # Sort the eigenvalues in descending order
70     sorted_eigenvalues, _ = torch.sort(eigenvalues, descending=True)
71
72     # Compute the cumulative explained variance
73     cumulative_explained_variances = torch.cumsum(sorted_eigenvalues, dim=0) / torch.sum(sorted_eigenvalues)
74
75     # Plotting
76     plt.figure(figsize=(10,6))
77     plt.plot(range(1, max_components+1), cumulative_explained_variances[:max_components].cpu().numpy(), marker='o')
78     plt.title("Cumulative Explained Variance vs. Number of Components")
79     plt.xlabel("Number of Components")
80     plt.ylabel("Cumulative Explained Variance")
81     plt.ylim(0, 1)
82     plt.grid(True)
83
84     save_plot("Cumulative_Variance", image_dir_path, plt)
```

```
85 plt.show()
86
87 def plot_clusters(data, max_clusters=10, image_dir_path=image_dir_path):
88     """
89     Display a line graph with points to visualize the sum of squared distances (SSD) for different cluster numbers.
90     """
91     # Compute the SSD for the range of clusters
92     ssd = [KMeans(n_clusters=k, n_init=10).fit(data).inertia_ for k in range(1, max_clusters+1)]
93
94     # Visualization
95     k_values = range(1, max_clusters+1)
96     plt.figure(figsize=(6, 6))
97     plt.plot(k_values, ssd, '-')
98     plt.scatter(k_values, ssd, color='b')
99     plt.title('Sum of Squared Distances vs Number of Clusters')
100    plt.xlabel('Number of clusters')
101    plt.ylabel('Sum of squared distances')
102    plt.grid(True)
103
104    save_plot("Cluster_Plot", image_dir_path, plt)
105
106    plt.show()
```

## Finding top 3 anomalous columns

```
In [39]: 1 # Combine training and validation data
2 df_scaled = scaler.fit_transform(df)
3
4 # Convert data to PyTorch tensors
5 dataset = TensorDataset(torch.tensor(df_scaled, dtype=torch.float32))
6 dataloader = DataLoader(dataset, batch_size=BATCH_SIZE, shuffle=False)
```

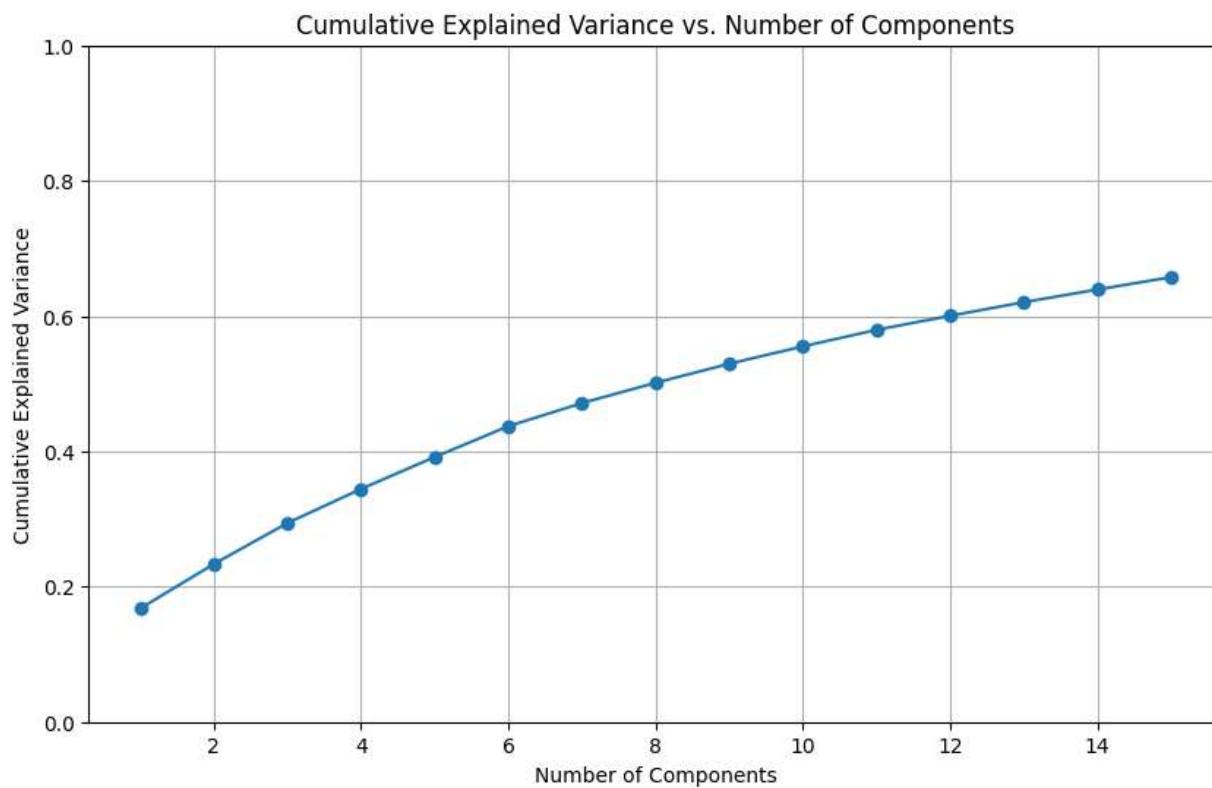
In [40]:

```
1 import torch
2
3 def calculate_reconstruction_and_top_3_values(model_name, model, dataloader, df_scaled, device):
4     reconstructed_batches = []
5     model.eval()
6     with torch.inference_mode():
7         for batch_data in tqdm(dataloader, desc=f'Processing {model_name}'):
8             inputs = batch_data[0].to(device)
9             reconstructed_batch = model(inputs)
10            reconstructed_batches.append(reconstructed_batch)
11
12    reconstructed = torch.cat(reconstructed_batches, dim=0).to(device)
13
14    # Reconstruction error calculation
15    reconstruction_error = ((reconstructed - torch.tensor(df_scaled, dtype=torch.float32).to(device)) ** 2).mean(dim=1).cpu()
16
17    # Top 3 values calculation
18    difference_percentage = 100 * (reconstructed - torch.tensor(df_scaled, dtype=torch.float32).to(device)) / (torch.tensor(
19        difference_df = pd.DataFrame(difference_percentage.cpu().numpy()))
20
21    # Ensure that difference_df is a tensor
22    difference_tensor = torch.tensor(difference_df.values).to(device)
23
24    # Compute the absolute value
25    abs_difference_tensor = torch.abs(difference_tensor)
26
27    # Determine the top 3 maximum absolute differences and their indices for each row
28    top_3_max_diff_values, indices = torch.topk(abs_difference_tensor, 3, dim=1)
29
30    # Convert the indices to column names
31    top_3_column_names_from_df = [[df.columns[i] for i in row] for row in indices.cpu().numpy()]
32
33    # Separate the values for each top 3 column
34    top_3_values = top_3_max_diff_values.cpu().numpy()
35
36    # Create a new DataFrame to store the top 3 max column names and values
37    top_3_max_diff_df = pd.DataFrame({
38        'Top_1_Column': [row[0] for row in top_3_column_names_from_df],
39        'Top_1_Value': top_3_values[:, 0],
40        'Top_2_Column': [row[1] for row in top_3_column_names_from_df],
41        'Top_2_Value': top_3_values[:, 1],
42        'Top_3_Column': [row[2] for row in top_3_column_names_from_df],
43        'Top_3_Value': top_3_values[:, 2],
44    })
45
46    return reconstruction_error, top_3_max_diff_df
```

## PCA Analysis

```
In [41]: 1 plot_cumulative_variance(df_scaled)
```

Image has been saved at Images\All\_Purchases\Cumulative\_Variance



```
In [42]: 1 dim_input = int(input("Please enter the number of dimensions: "))
```

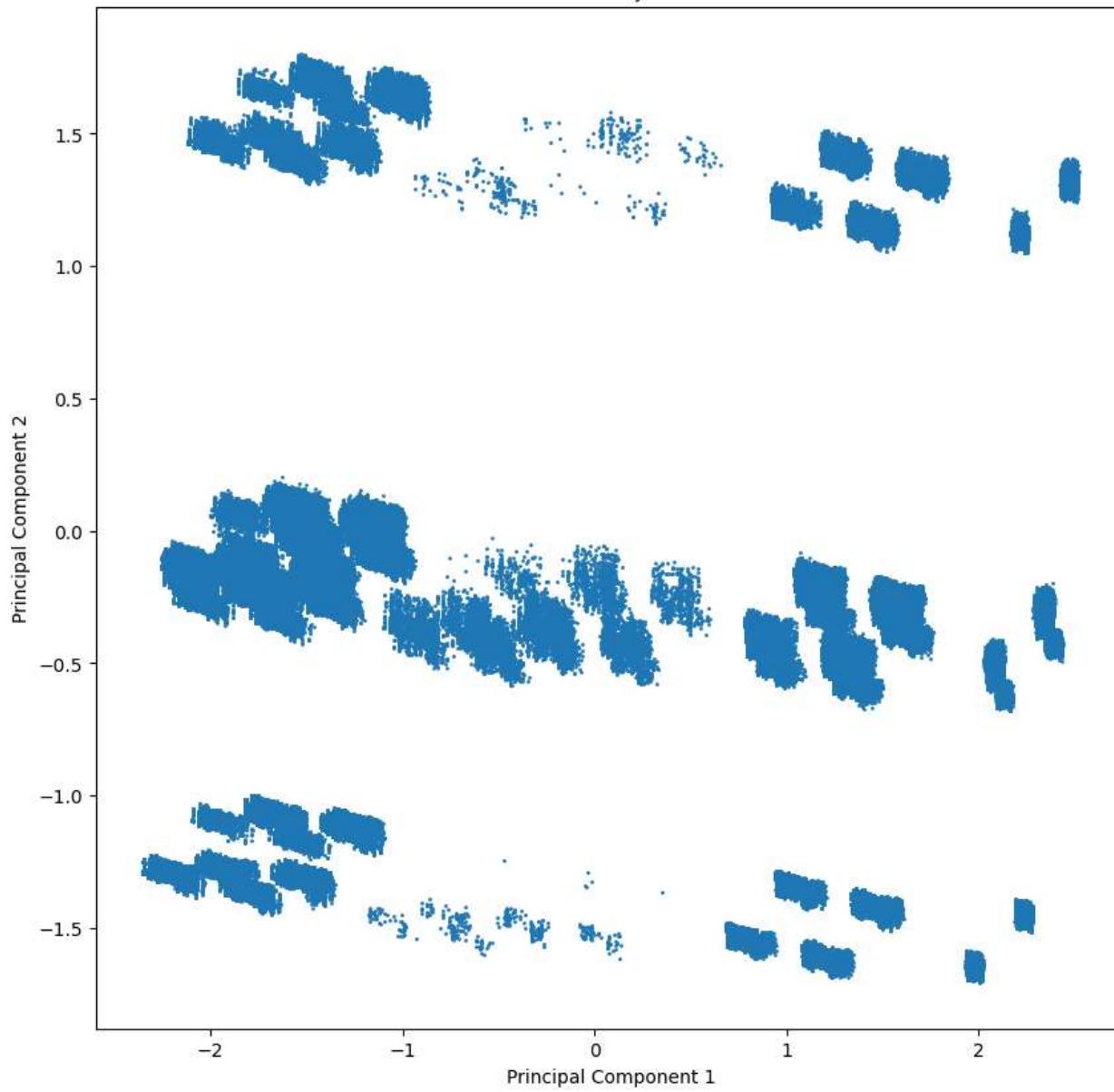
Please enter the number of dimensions: 9

In [43]:

```
1 # PCA Analysis
2 import warnings
3
4 with warnings.catch_warnings():
5     warnings.simplefilter("ignore")
6     transformed_data = apply_pca_optimized(df_scaled, n_components=dim_input)
```

Image has been saved at Images\All\_Purchases\PCA\_Visualization

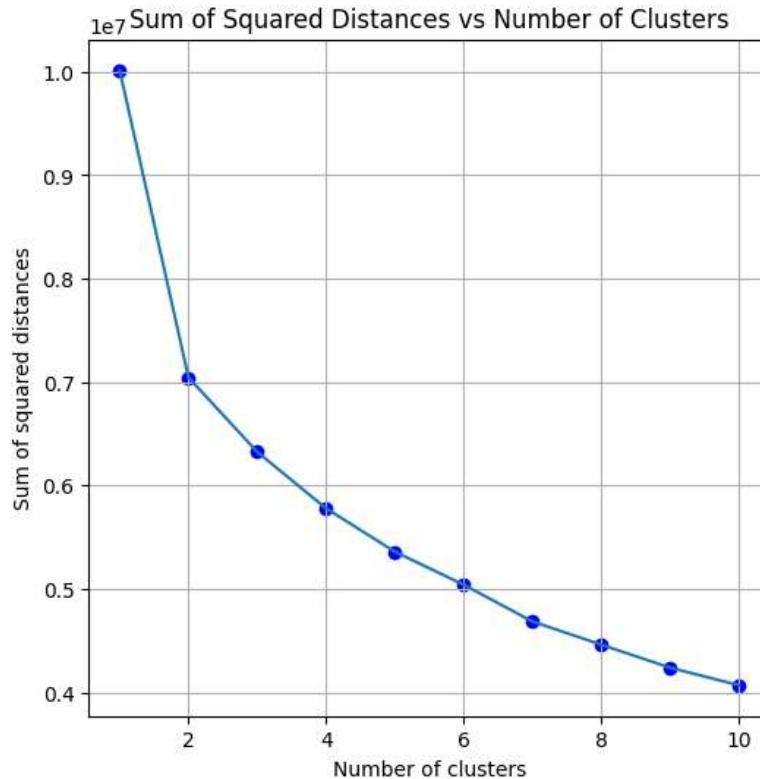
2D PCA Projection



Total Explained Variance (9 components): 52.95%

```
In [44]: 1 # Visualize Cluster SSD  
2 plot_clusters(transformed_data)
```

Image has been saved at Images\All\_Purchases\Cluster\_Plot



```
In [45]: 1 cluster_input = int(input("Please enter the number of clusters: "))
```

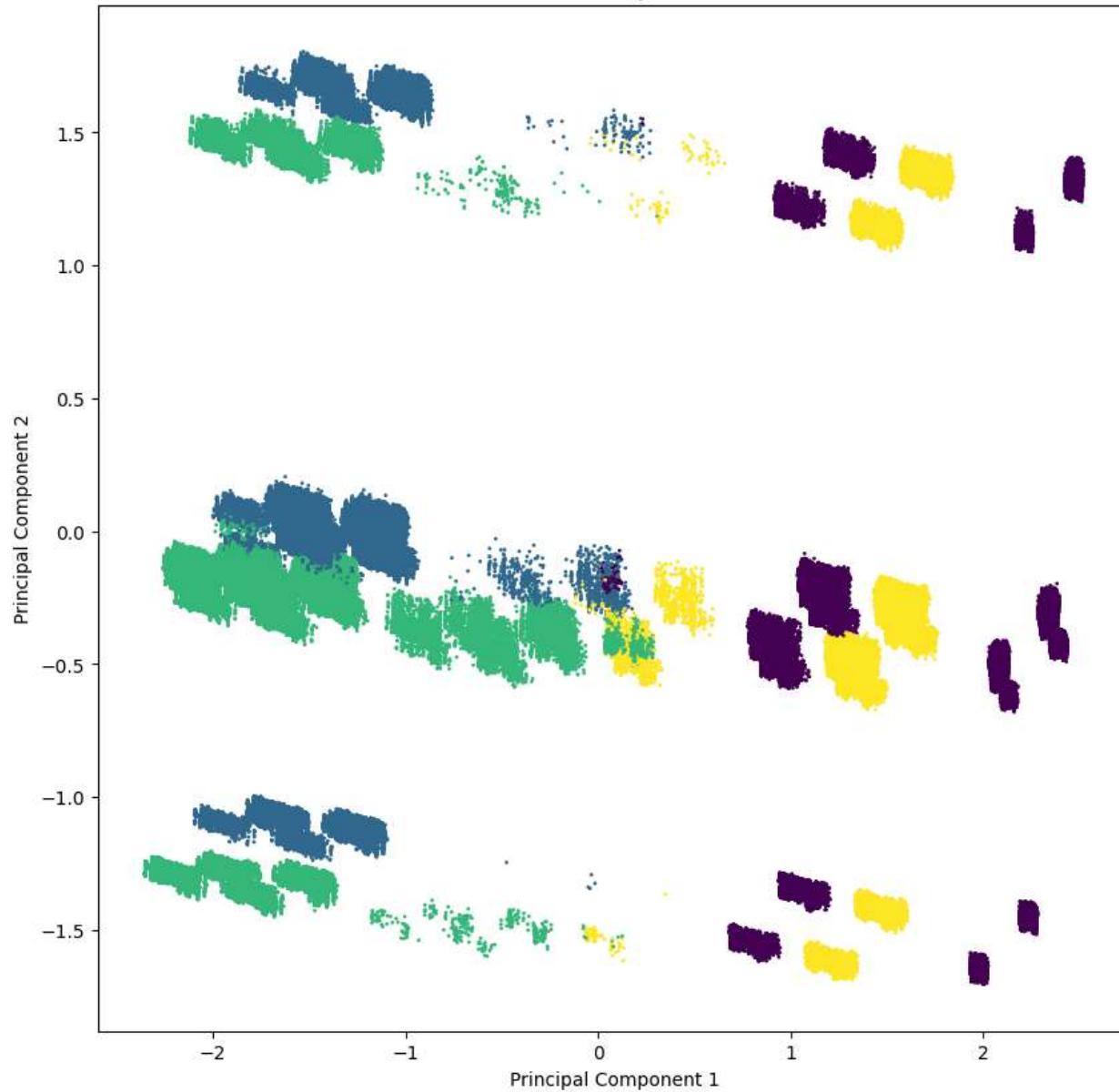
Please enter the number of clusters: 4

```
In [46]: 1 kmeans = KMeans(n_clusters=cluster_input, n_init=10)  
2 kmeans.fit(transformed_data)  
3 labels = kmeans.labels_  
4 df_labels = pd.DataFrame(labels, columns=['PCA_Labels'])
```

```
In [47]: 1 visualize_pca(transformed_data, labels=labels)
```

Image has been saved at Images\All\_Purchases\PCA\_Visualization\_Labels

2D PCA Projection



```
In [48]: 1 df_transformed_data = pd.DataFrame(transformed_data)
```

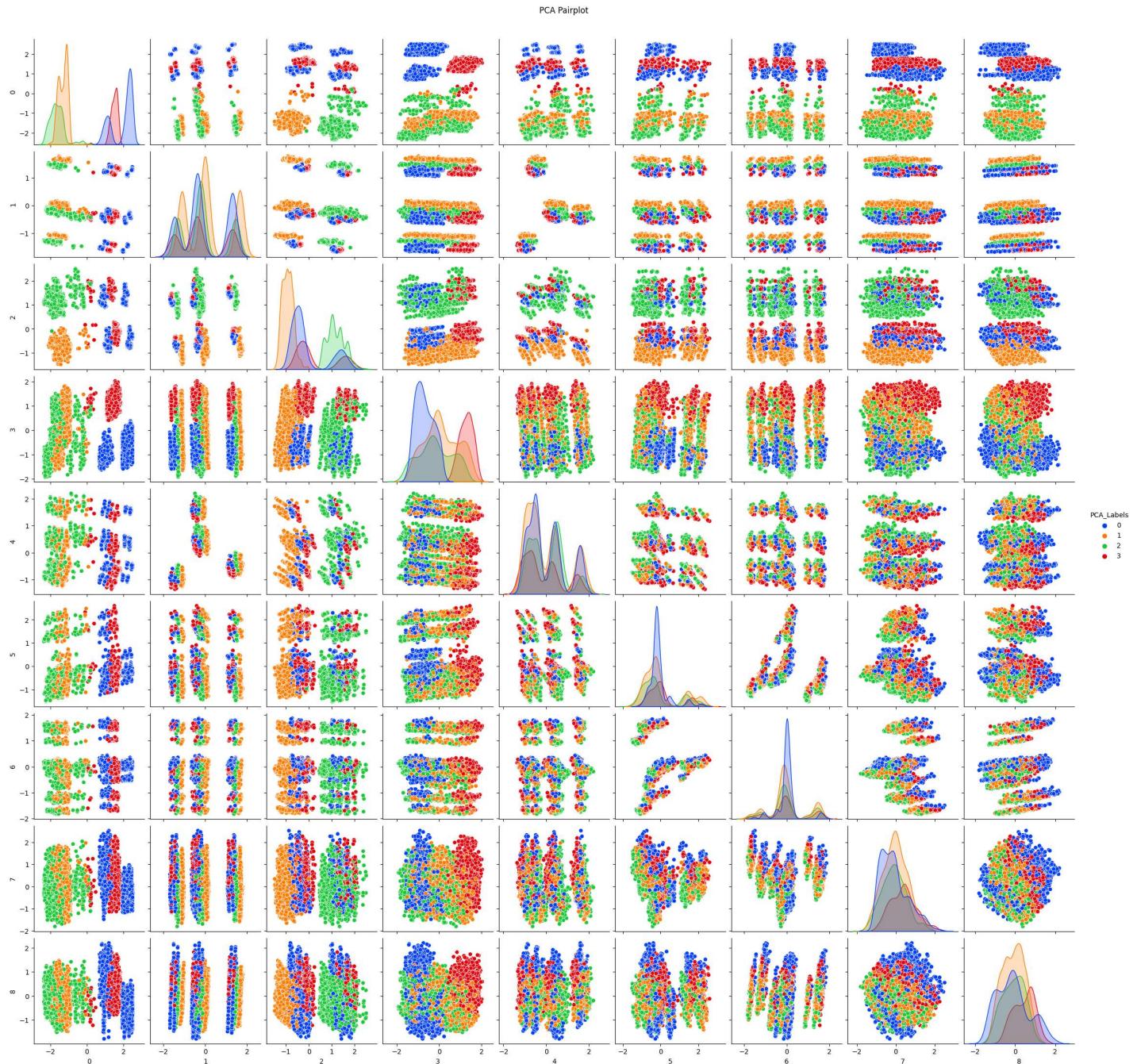
```
2 df_transformed_data['PCA_Labels'] = labels
```

```
3 sampled_data = df_transformed_data.sample(n=10000)
```

In [49]:

```
1 import seaborn as sns
2
3 # Create the pairplot
4 with warnings.catch_warnings():
5     warnings.simplefilter("ignore")
6     plot = sns.pairplot(sampled_data, hue='PCA_Labels', palette="bright")
7
8 # Adjusting space at the top for the title
9 plt.subplots_adjust(top=0.95)
10 plt.suptitle("PCA Pairplot")
11
12 # Call the save_plot function to save the image, with the correct file extension
13 save_plot("PCA_Pairplot", image_dir_path, plot.fig)
```

Image has been saved at Images\All\_Purchases\PCA\_Pairplot



## Finding anomalies

In [50]:

```
1 from tqdm.notebook import tqdm
2 import os
3 import numpy as np
4 import pandas as pd
5 from datetime import datetime, timedelta
6
7 def detect_anomalies(reconstruction_error, threshold):
8     anomaly_threshold = np.percentile(reconstruction_error, threshold)
9     return reconstruction_error > anomaly_threshold
10
11 def create_dataframe(reconstruction_error, df_orig, df_top_3, df_labels):
12     scores_df = pd.DataFrame(reconstruction_error, columns=['anomaly_score'])
13     df_score = pd.concat([scores_df, df_orig, df_top_3, df_labels], axis=1)
14     return df_score
15
16 def save_anomalies(all_anomalies, model_name, filename, result_dir_path=None):
17     if result_dir_path is None:
18         raise ValueError("result_dir_path must be provided.")
19     dir_path = os.path.join(result_dir_path, model_name)
20     os.makedirs(dir_path, exist_ok=True)
21     all_anomalies.to_csv(os.path.join(dir_path, filename), index=False)
22
23 def filter_recent_anomalies(all_anomalies):
24     cutoff_date = datetime.now() - timedelta(days=30)
25     all_anomalies['sls_d'] = pd.to_datetime(all_anomalies['sls_d'])
26     df_filtered = all_anomalies[all_anomalies['sls_d'] >= cutoff_date]
27     return df_filtered.sort_values(by=['anomaly_score', 'sls_d', 'start_tm'], ascending=[False, False, True])
28
29 def detect_and_save_anomalies(model_name, model, dataloader, df_top_3, labels, df_orig, threshold=99.9, result_dir_path=None):
30     if reconstruction_error is None:
31         raise ValueError("reconstruction_error must be provided.")
32     if result_dir_path is None:
33         raise ValueError("result_dir_path must be provided.")
34     print(f"{model_name} has estimated outputs")
35     anomalies = detect_anomalies(reconstruction_error, threshold)
36     print(f"{model_name} reconstruction error calculated")
37     df_score = create_dataframe(reconstruction_error, df_orig, df_top_3, pd.DataFrame(labels, columns=['PCA_Labels'])) # Add
38     all_anomalies = df_score[df_score.index.isin(df_orig[anomalies].index)].sort_values(by='anomaly_score', ascending=False)
39     print(f"{model_name} anomaly dataframe created")
40     save_anomalies(all_anomalies, model_name, "All_Anomalies.csv", result_dir_path)
41     print(f"{model_name} anomaly csv saved")
42     df_filtered = filter_recent_anomalies(all_anomalies)
43     # No index resetting here
44     save_anomalies(df_filtered, model_name, "Recent_Anomalies.csv", result_dir_path)
45     print(f"{model_name} recent anomaly csv saved\n")
```

## Export to CSV

In [51]:

```
1 for key, value in models.items():
2     reconstruction_error, top_3_max_diff_df = calculate_reconstruction_and_top_3_values(key, value, dataloader, df_scaled,
3         detect_and_save_anomalies(key, value, dataloader, top_3_max_diff_df, labels, df_orig, 99.9, result_dir_path=result_dir_
```

A Jupyter widget could not be displayed because the widget state could not be found. This could happen if the kernel storing the widget is no longer available, or if the widget state was not saved in the notebook. You may be able to create the widget by running the appropriate cells.

```
model_1 has estimated outputs
model_1 reconstruction error calculated
model_1 anomaly dataframe created
model_1 anomaly csv saved
model_1 recent anomaly csv saved
```

A Jupyter widget could not be displayed because the widget state could not be found. This could happen if the kernel storing the widget is no longer available, or if the widget state was not saved in the notebook. You may be able to create the widget by running the appropriate cells.

```
model_2 has estimated outputs
model_2 reconstruction error calculated
model_2 anomaly dataframe created
model_2 anomaly csv saved
model_2 recent anomaly csv saved
```

A Jupyter widget could not be displayed because the widget state could not be found. This could happen if the kernel storing the widget is no longer available, or if the widget state was not saved in the notebook. You may be able to create the widget by running the appropriate cells.

```
model_3 has estimated outputs
model_3 reconstruction error calculated
model_3 anomaly dataframe created
model_3 anomaly csv saved
model_3 recent anomaly csv saved
```

A Jupyter widget could not be displayed because the widget state could not be found. This could happen if the kernel storing the widget is no longer available, or if the widget state was not saved in the notebook. You may be able to create the widget by running the appropriate cells.

```
model_4 has estimated outputs
model_4 reconstruction error calculated
model_4 anomaly dataframe created
model_4 anomaly csv saved
model_4 recent anomaly csv saved
```

In [ ]:

1