





```
In [40]: train_set.iloc[:, :405].describe(include="all")
```

	V338	V337	V338	V339	id_01	id_02	id_03	id_04	id_05
count	12268.0	12268.000000	12268.000000	12268.000000	1.645900e+04	16069.000000	7456.000000	7456.000000	15605.000000
mean	inf	15.742602	21.362671	16.464481	-inf	133358.390625	0.107544	-0.070801	1.869648
std	inf	91.064587	131.079559	92.732029	1.403906e+01	123074.875000	0.707195	0.726662	5.230469
min	0.0	0.000000	0.000000	0.000000	-1.000000e+02	1120.000000	-8.000000	-13.000000	-80.000000
25%	0.0	0.000000	0.000000	0.000000	-5.000000e+00	52145.000000	0.000000	0.000000	0.000000
50%	0.0	0.000000	0.000000	0.000000	-5.000000e+00	97865.000000	0.000000	0.000000	0.000000
75%	0.0	0.000000	0.000000	0.000000	0.000000e+00	171732.000000	0.000000	0.000000	1.000000
max	1300.0	1700.000000	3000.000000	1700.000000	0.000000e+00	998295.000000	9.000000	0.000000	52.000000

```
In [41]: train_set.iloc[:, 405:420].describe(include="all")
```

	id_12	id_13	id_14	id_15	id_16	id_17	id_18	id_19	id_20	id_21	id_22	id_23	id_24
count	16456	12762.0	11894.0	10069	15137	15921.0	4853.0	15917.0	15912.0	604.0	608.0	606	559.0
unique	2	36.0	21.0	3	2	49.0	12.0	349.0	164.0	136.0	6.0	3	5.0
top	NotFound	52.0	-300.0	Found	NotFound	166.0	15.0	410.0	567.0	252.0	14.0	IP_PROXY:TRANSPARENT	15.0
freq	14433	6997.0	6565.0	7786	7681	11828.0	2342.0	1693.0	1828.0	200.0	559.0	329	265.0

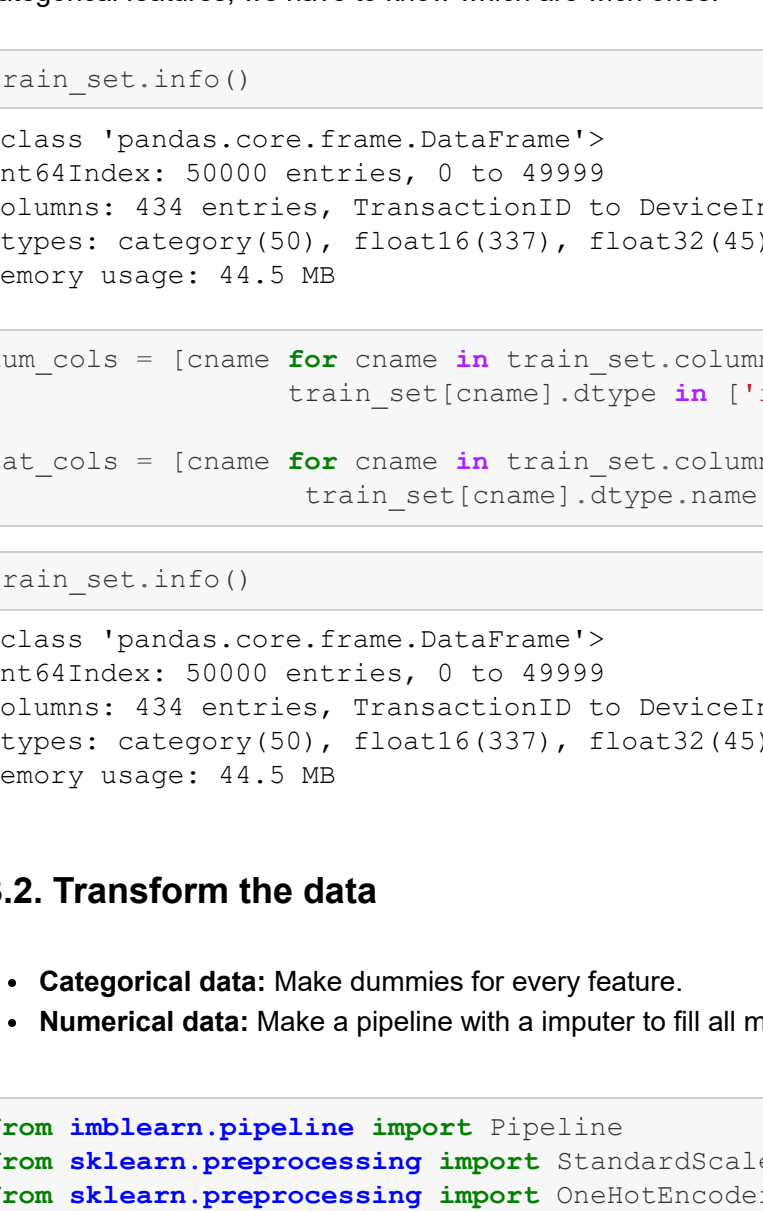
```
In [42]: train_set.iloc[:, 420:435].describe(include="all")
```

	id_27	id_28	id_29	id_30	id_31	id_32	id_33	id_34	id_35	id_36	id_37	id_38	DeviceType	DeviceInfo
count	608	16069	16069	11673	16048	11674.0	10821	11640	16069	16069	16069	16069	16068	14351
unique	2	2	2	63	73	4.0	89	3	2	2	2	2	2	585
top	Found	Found	Found	Windows	chrome	62.0	24.0	1920x1080	match_status_2	T	F	T	desktop	Windows
freq	607	9059	8750	3344	4956	8129.0	2576	9104	11676	14703	12671	13423	10887	6224

### 3. Data cleaning

```
In [43]: counts = train_set['isFraud'].value_counts().values
ax = sns.barplot(0, counts)
ax.set(title="Class distribution w.r.t target variables", xlabel = 'Target Class Count', ylabel='Count')

Out[43]: (Text(0, 0.5, 'Count'),
Text(0.5, 0, 'Target Class Count'),
Text(0.5, 1.0, 'Class distribution w.r.t target variables'))
```



```
In [44]: train_set.isnull().sum() / train_set.isnull().sum() != 0]
```

```
Out[44]: card2      697
card3           3
card4           6
card5          235
card6           3
      ...
id_36      33931
id_37      33931
id_38      33931
DeviceType   33932
DeviceInfo   35649
Length: 359, dtype: int64
```

#### 3.1. split-up categorical and numerical features

**note:** We see that there is almost missing data in e every feature. Before we can scale and impute our numerical and one-hot-encode our categorical features, we have to know which are wich ones.

```
In [45]: train_set.info()
```

```
<class 'pandas.core.frame.DataFrame'>
Int64Index: 50000 entries, 0 to 49999
Columns: 434 entries, TransactionID to DeviceInfo
dtypes: category(50), float16(337), float32(45), int32(2)
memory usage: 44.5 MB
```

```
In [46]: num_cols = [cname for cname in train_set.columns if
train_set[cname].dtype in ['int32', 'float16', 'float32']]

cat_cols = [cname for cname in train_set.columns if
train_set[cname].dtype.name == "category"]
```

```
In [47]: train_set.info()
```

```
<class 'pandas.core.frame.DataFrame'>
Int64Index: 50000 entries, 0 to 49999
Columns: 434 entries, TransactionID to DeviceInfo
dtypes: category(50), float16(337), float32(45), int32(2)
memory usage: 44.5 MB
```

#### 3.2. Transform the data

- Categorical data: Make a pipeline for every feature.
- Numerical data: Make a pipeline with a imputer to fill all missing data. We'll also scale the data.

```
In [48]: from imblearn.pipeline import Pipeline
from sklearn.preprocessing import StandardScaler
from sklearn.preprocessing import OneHotEncoder
from sklearn.impute import SimpleImputer
from sklearn.compose import ColumnTransformer
```

```
trainset = train_set.copy()
testset = test_set.copy()
```

```
#Preprocessing for categorical data
for col in cat_cols:
    if col != 'isFraud':
        dummies = pd.get_dummies(trainset[col], dummy_na=False, prefix=col)
        trainset = pd.concat([trainset, dummies], axis=1)
        trainset.drop(columns=[col], inplace=True)

    testset = pd.concat([testset, dummies], axis=1)
    testset.drop(columns=[col], inplace=True)
```

```
#Preprocessng for numerical data
numeric_transformer = Pipeline(steps=[
    ('imputer', SimpleImputer(strategy='mean')),
    ('scaler', StandardScaler())])

trainset[num_cols] = numeric_transformer.fit_transform(trainset[num_cols])
testset[num_cols] = numeric_transformer.fit_transform(testset[num_cols])
```

```
In [49]: trainset.isnull().sum() / trainset.isnull().sum() != 0]
```

```
Out[49]: Series(1), dtype: int64
```

```
In [50]: trainset.shape
```

```
Out[50]: (50000, 8689)
```

```
In [51]: testset.shape
```

```
Out[51]: (50000, 8688)
```

### 4. Building auto-encoder

#### 4.1. Train - test - val - split

```
In [52]: X_train, X_test = train_test_split(trainset, test_size=0.2, random_state=1234)
```

```
Y_test = X_test['isFraud']
X_test = X_test.drop(['isFraud'], 1)

X_train = X_train[X_train['isFraud'] == 0]
X_train = X_train.drop(['isFraud'], 1)

X_val, X_test, y_val, y_test = train_test_split(X_train, Y_test, test_size=0.4, random_state=123)
```

```
X_train = X_train.values
X_val = X_val.values
X_test = X_test.values
```

We'll check the amount of observations that are flagged as fraud, so we know how much observations we are looking for.

```
In [53]: y_val.value_counts()
```

```
Out[53]: 0      5846
1         154
Name: isFraud, dtype: int64
```

#### 4.2. Building the auto-encoder model

```
In [54]: # the number of features is the input shape
input_dim = X_train.shape[1]
```

```
# rate
rate = 0.2

# encoder model
# Input Layer
input_layer = Input(shape=(input_dim, ))

# encoder layers
encoder = Dense(128, activation="relu", activity_regularizer=regularizers.l1(10e-5))(input_layer)
BatchNormalization()(encoder)
encoder = Dense(64, activation="relu")(encoder)
Dropout(rate)(encoder)
BatchNormalization()(encoder)
encoder = Dense(input_dim, activation="sigmoid")(decoder)

# decoder model
decoder = Dense(64, activation = "relu" )(encoder)
Dropout(rate)(decoder)
BatchNormalization()(decoder)
decoder = Dense(128, activation="relu" )(encoder)
Dropout(rate)(decoder)
BatchNormalization()(decoder)
decoder = Dense(input_dim, activation="sigmoid" )(decoder)

# build auto-encoder model
autoencoder = Model(inputs = input_layer, outputs = decoder)
```

#### 4.3. Compiling and fitting the model

```
In [55]: #compile and fit
from tensorflow.python.keras.callbacks import ModelCheckpoint
from tensorflow.python.keras.callbacks import EarlyStopping
```

```
nb_epoch = 100
batch_size = 29
early_stopping = EarlyStopping(patience=10)
autoencoder.compile(optimizer='adam',
                    loss='mean_squared_error',
                    metrics=['accuracy'])

checkpointer = ModelCheckpoint(filepath="model.h5",
                               verbose=0,
                               save_best_only=True)
```

```
history = autoencoder.fit(X_train, X_train,
                          epochs = nb_epoch,
                          batch_size = batch_size,
                          shuffle=True,
                          validation_split=0.2,
                          verbose=1,
                          callbacks=[checkpointer, early_stopping])
```

```
Epoch 1/100
1074/1074 [=====] - 32s 29ms/step - loss: 0.0613 - accuracy: 4.3939e-04 - val_loss: 0.0357 - val_accuracy: 0.0037
Epoch 2/100
1074/1074 [=====] - 30s 28ms/step - loss: 0.0341 - accuracy: 0.0056 - val_loss: 0.0339 - val_accuracy: 0.0030
Epoch 3/100
1074/1074 [=====] - 24s 22ms/step - loss: 0.0339 - accuracy: 0.0084 - val_loss: 0.0348 - val_accuracy: 0.0267
Epoch 4/100
1074/1074 [=====] - 25s 23ms/step - loss: 0.0327 - accuracy: 0.0460 - val_loss: 0.0345 - val_accuracy: 0.0123
Epoch 5/100
1074/1074 [=====] - 27s 25ms/step - loss: 0.0320 - accuracy: 0.0364 - val_loss: 0.0344 - val_accuracy: 0.0444
Epoch 6/100
1074/1074 [=====] - 25s 23ms/step - loss: 0.0328 - accuracy: 0.0381 - val_loss: 0.0348 - val_accuracy: 0.0460
Epoch 7/100
1074/1074 [=====] - 30s 28ms/step - loss: 0.0325 - accuracy: 0.0351 - val_loss: 0.0342 - val_accuracy: 0.0325
Epoch 8/100
1074/1074 [=====] - 33s 31ms/step - loss: 0.0313 - accuracy: 0.0302 - val_loss: 0.0341 - val_accuracy: 0.0288
Epoch 9/100
1074/1074 [=====] - 30s 28ms/step - loss: 0.0319 - accuracy: 0.0345 - val_loss: 0.0340 - val_accuracy: 0.0374
Epoch 10/100
1074/1074 [=====] - 33s 30ms/step - loss: 0.0319 - accuracy: 0.0278 - val_loss: 0.0340 - val_accuracy: 0.0324
Epoch 11/100
1074/1074 [=====] - 34s 32ms/step - loss: 0.0325 - accuracy: 0.0291 - val_loss: 0.0339 - val_accuracy: 0.0320
Epoch 12/100
1074/1074 [=====] - 31s 29ms/step - loss: 0.0341 - accuracy: 0.0286 - val_loss: 0.0339 - val_accuracy: 0.0380
Epoch 13/100
1074/1074 [=====] - 29s 27ms/step - loss: 0.0321 - accuracy: 0.0350 - val_loss: 0.0345 - val_accuracy: 0.0619
Epoch 14/100
1074/1074 [=====] - 27s 25ms/step - loss: 0.0334 - accuracy: 0.0370 - val_loss: 0.0344 - val_accuracy: 0.0378
Epoch 15/100
1074/1074 [=====] - 28s 26ms/step - loss: 0.0341 - accuracy: 0.0350 - val_loss: 0.0338 - val_accuracy: 0.0560
Epoch 16/100
1074/1074 [=====] - 28s 26ms/step - loss: 0.0318 - accuracy: 0.0415 - val_loss: 0.0338 - val_accuracy: 0.0348
Epoch 17/100
1074/1074 [=====] - 30s 28ms/step - loss: 0.0319 - accuracy: 0.0357 - val_loss: 0.0338 - val_accuracy: 0.0424
Epoch 18/100
1074/1074 [=====] - 30s 28ms/step - loss: 0.0307 - accuracy: 0.0380 - val_loss: 0.0338 - val_accuracy: 0.0378
Epoch 19/100
1074/1074 [=====] - 30s 28ms/step - loss: 0.0323 - accuracy: 0.0394 - val_loss: 0.0337 - val_accuracy: 0.0600
Epoch 20/100
1074/1074 [=====] - 30s 28ms/step - loss: 0.0320 - accuracy: 0.0408 - val_loss: 0.0337 - val_accuracy: 0.0487
Epoch 21/100
1074/1074 [=====] - 34s 32ms/step - loss: 0.0323 - accuracy: 0.0399 - val_loss: 0.0337 - val_accuracy: 0.0597
Epoch 22/100
1074/1074 [=====] - 37s 34ms/step - loss: 0.0338 - accuracy: 0.0425 - val_loss: 0.0336 - val_accuracy: 0.0528
Epoch 23/100
1074/1074 [=====] - 35s 32ms/step - loss: 0.0322 - accuracy: 0.0478 - val_loss: 0.0336 - val_accuracy: 0.0533
Epoch 24/100
1074/1074 [=====] - 29s 27ms/step - loss: 0.0310 - accuracy: 0.0482 - val_loss: 0.0336 - val_accuracy: 0.0446
Epoch 25/100
1074/1074 [=====] - 29s 27ms/step - loss: 0.0322 - accuracy: 0.0447 - val_loss: 0.0336 - val_accuracy: 0.0438
Epoch 26/100
1074/1074 [=====] - 29s 27ms/step - loss: 0.0321 - accuracy: 0.0480 - val_loss: 0.0336 - val_accuracy: 0.0551
Epoch 27/100
1074/1074 [=====] - 33s 30ms/step - loss: 0.0327 - accuracy: 0.0489 - val_loss: 0.0336 - val_accuracy: 0.0555
Epoch 28/100
1074/1074 [=====] - 32s 30ms/step - loss: 0.0319 - accuracy: 0.0486 - val_loss: 0.0336 - val_accuracy: 0.0672
Epoch 29/100
1074/1074 [=====] - 33s 30ms/step - loss: 0.0326 - accuracy: 0.0553 - val_loss: 0.0336 - val_accuracy: 0.0645
Epoch 30/100
1074/1074 [=====] - 33s 31ms/step - loss: 0.0321 - accuracy: 0.0481 - val_loss: 0.0336 - val_accuracy: 0.0545
Epoch 31/100
1074/1074 [=====] - 31s 29ms/step - loss: 0.0308 - accuracy: 0.0535 - val_loss: 0.0336 - val_accuracy: 0.0629
Epoch 32/100
1074/1074 [=====] - 30s 28ms/step - loss: 0.0310 - accuracy: 0.0589 - val_loss: 0.0336 - val_accuracy: 0.0704
Epoch 33/100
1074/1074 [=====] - 28s 26ms/step - loss: 0.0322 - accuracy: 0.0573 - val_loss: 0.0336 - val_accuracy: 0.0697
Epoch 34/100
1074/1074 [=====] - 28s 26ms/step - loss: 0.0313 - accuracy: 0.0560 - val_loss: 0.0336 - val_accuracy: 0.0626
Epoch 35/100
1074/1074 [=====] - 28s 26ms/step - loss: 0.0317 - accuracy: 0.0653 - val_loss: 0.0336 - val_accuracy: 0.0613
Epoch 36/100
1074/1074 [=====] - 33s 31ms/step - loss: 0.0323 - accuracy: 0.0591 - val_loss: 0.0336 - val_accuracy: 0.0764
Epoch 37/100
1074/1074 [=====] - 37s 34ms/step - loss: 0.0321 - accuracy: 0.0584 - val_loss: 0.0336 - val_accuracy: 0.0533
Epoch 38/100
1074/1074 [=====] - 34s 31ms/step - loss: 0.0312 - accuracy: 0.0540 - val_loss: 0.0336 - val_accuracy: 0.0749
Epoch 39/100
1074/1074 [=====] - 32s 30ms/step - loss: 0.0316 - accuracy: 0.0611 - val_loss: 0.0336 - val_accuracy: 0.0906
Epoch 40/100
1074/1074 [=====] - 24s 22ms/step - loss: 0.0336 - accuracy: 0.0589 - val_loss: 0.0336 - val_accuracy: 0.0678
Epoch 41/100
1074/1074 [=====] - 35s 33ms/step - loss: 0.0320 - accuracy: 0.0583 - val_loss: 0.0336 - val_accuracy: 0.0681
Epoch 42/100
1074/1074 [=====] - 36s 33ms/step - loss: 0.0319 - accuracy: 0.0558 - val_loss: 0.0336 - val_accuracy: 0.0746
Epoch 43/100
1074/1074 [=====] - 30s 28ms/step - loss: 0.0325 - accuracy: 0.0622 - val_loss: 0.0336 - val_accuracy: 0.0825
Epoch 44/100
1074/1074 [=====] - 35s 33ms/step - loss: 0.0301 - accuracy: 0.0690 - val_loss: 0.0336 - val_accuracy: 0.0858
Epoch 45/100
1074/1074 [=====] - 36s 33ms/step - loss: 0.0327 - accuracy: 0.0657 - val_loss: 0.0336 - val_accuracy: 0.0897
Epoch 46/100
1074/1074 [=====] - 34s 31ms/step - loss: 0.0335 - accuracy: 0.0861 - val_loss: 0.0336 - val_accuracy: 0.0754
Epoch 47/100
1074/1074 [=====] - 37s 34ms/step - loss: 0.0317 - accuracy: 0.0877 - val_loss: 0.0336 - val_accuracy: 0.0764
Epoch 48/100
1074/1074 [=====] - 31s 29ms/step - loss: 0.0320 - accuracy: 0.0812 - val_loss: 0.0336 - val_accuracy: 0.0782
Epoch 49/100
1074/1074 [=====] - 21s 20ms/step - loss: 0.0311 - accuracy: 0.0824 - val_loss: 0.0336 - val_accuracy: 0.0746
Epoch 50/100
1074/1074 [=====] - 10s 9ms/step - loss: 0.0312 - accuracy: 0.0720 - val_loss: 0.0336 - val_accuracy: 0.0845
Epoch 51/100
1074/1074 [=====] - 8s 7ms/step - loss: 0.0325 - accuracy: 0.0759 - val_loss: 0.0336 - val_accuracy: 0.0929
Epoch 52/100
1074/1074 [=====] - 8s 7ms/step - loss: 0.0324 - accuracy: 0.0791 - val_loss: 0.0336 - val_accuracy: 0.0948
Epoch 53/100
1074/1074 [=====] - 8s 7ms/step - loss: 0.0307 - accuracy: 0.0764 - val_loss: 0.0336 - val_accuracy: 0.0827
Epoch 54/100
1074/1074 [=====] - 8s 7ms/step - loss: 0.0307 - accuracy: 0.0708 - val_loss: 0.0336 - val_accuracy: 0.0897
Epoch 55/100
1074/1074 [=====] - 8s 7ms/step - loss: 0.0328 - accuracy: 0.0896 - val_loss: 0.0336 - val_accuracy: 0.1127
Epoch 56/100
1074/1074 [=====] - 8s 7ms/step - loss: 0.0309 - accuracy: 0.0882 - val_loss: 0.0336 - val_accuracy: 0.1037
Epoch 57/100
1074/1074 [=====] - 8s 7ms/step - loss: 0.0326 - accuracy: 0.0908 - val_loss: 0.0336 - val_accuracy: 0.1150
Epoch 58/100
1074/1074 [=====] - 8s 7ms/step - loss: 0.0304 - accuracy: 0.0870 - val_loss: 0.0336 - val_accuracy: 0.1021
Epoch 59/100
1074/1074 [=====] - 8s 7ms/step - loss: 0.0310 - accuracy: 0.0845 - val_loss: 0.0336 - val_accuracy: 0.0999
Epoch 60/100
1074/1074 [=====] - 8s 7ms/step - loss: 0.0302 - accuracy: 0.0896 - val_loss: 0.0336 - val_accuracy: 0.1010
Epoch 61/100
1074/1074 [=====] - 8s 7ms/step - loss: 0.0319 - accuracy: 0.0959 - val_loss: 0.0336 - val_accuracy: 0.0832
Epoch 62/100
1074/1074 [=====] - 8s 7ms/step - loss: 0.0308 - accuracy: 0.0791 - val_loss: 0.0336 - val_accuracy: 0.0829
Epoch 63/100
1074/1074 [=====] - 8s 7ms/step - loss: 0.0323 - accuracy: 0.0841 - val_loss: 0.0336 - val_accuracy: 0.0972
Epoch 64/100
1074/1074 [=====] - 8s 7ms/step - loss: 0.0319 - accuracy: 0.0823 - val_loss: 0.0336 - val_accuracy: 0.1083
Epoch 65/100
1074/1074 [=====] - 8s 7ms/step - loss: 0.0328 - accuracy: 0.0832 - val_loss: 0.0336 - val_accuracy: 0.1021
Epoch 66/100
1074/1074 [=====] - 8s 7ms/step - loss: 0.0325 - accuracy: 0.0953 - val_loss: 0.0336 - val_accuracy: 0.1019
Epoch 67/100
1074/1074 [=====] - 8s 7ms/step - loss: 0.0331 - accuracy: 0.1001 - val_loss: 0.0336 - val_accuracy: 0.1152
Epoch 68/100
1074/1074 [=====] - 8s 7ms/step - loss: 0.0324 - accuracy: 0.1010 - val_loss: 0.0336 - val_accuracy: 0.1132
Epoch 69/100
1074/1074 [=====] - 8s 7ms/step - loss: 0.0309 - accuracy: 0.0965 - val_loss: 0.0336 - val_accuracy: 0.1107
Epoch 70/100
1074/1074 [=====] - 8s 7ms/step - loss: 0.0305 - accuracy: 0.1023 - val_loss: 0.0336 - val_accuracy: 0.1075
Epoch 71/100
1074/1074 [=====] - 8s 7ms/step - loss: 0.0317 - accuracy: 0.1027 - val_loss: 0.0336 - val_accuracy: 0.1236
Epoch 72/100
1074/1074 [=====] - 8s 7ms/step - loss: 0.0308 - accuracy: 0.1029 - val_loss: 0.0336 - val_accuracy: 0.1141
```

#### 4.4. Predicting/ Reconstructing credit card transactions trainset and testset

```
In [56]: predictions_train = autoencoder.predict(X_train)
predictions_val = autoencoder.predict(X_val)
```

#### 4.5. Taking a look at the MSE

```
In [57]: from sklearn.metrics import mean_squared_error

mse_val = mean_squared_error(X_val, predictions_val)
mse_train = mean_squared_error(X_train, predictions_train)
print(mse_train)
print(mse_val)

0.03191511570509188
0.03689375289796546
```

```
In [58]: error_train = np.mean(np.power(X_train - predictions_train, 2), axis=1)
limit = np.percentile(error_train, 97.5)
print(limit)
```

```
0.17885801674083593
```

```
In [59]: error_val = np.mean(np.power(X_val - predictions_val, 2), axis=1)
fraud = error_val > limit
fraud.sum()
```

```
Out[59]: 175
```

```
In [60]: error_val_pd = pd.DataFrame(error_val)
error_val_pd.columns = ['error']
error_val_pd['id'] = range(len(error_val_pd))
error_val_pd.info()
```

```
<class 'pandas.core.frame.DataFrame'>
RangeIndex: 6000 entries, 0 to 5999
Data columns (total 2 columns):
 # Column Non-Null Count Dtype
 0 error 6000 non-null float64
 1 id 6000 non-null int32
dtypes: float64(1), int32(1)
memory usage: 10.4 KB
```

```
In [61]: from sklearn.metrics import confusion_matrix, precision_recall_curve

confusion_mat = confusion_matrix(y_val, fraud)
print(confusion_mat)
```

```
[1596 150]
[129 25]
```

```
In [62]: from sklearn.metrics import classification_report
print(classification_report(y_val, fraud))
```

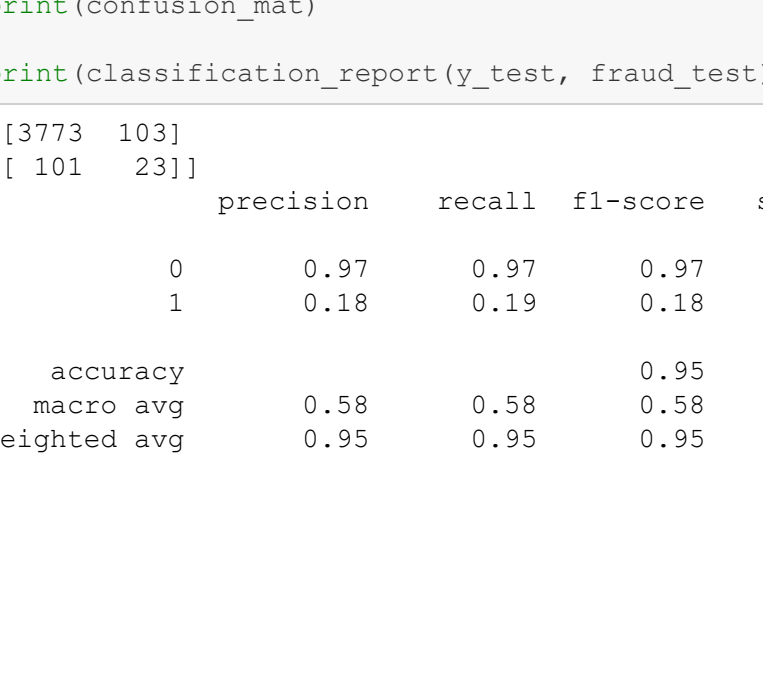
```
precision recall f1-score support

0 0.98 0.97 0.98 5846
1 0.14 0.16 0.15 154

accuracy 0.95 6000
macro avg 0.56 0.57 0.56 6000
weighted avg 0.96 0.95 0.95 6000
```

```
In [63]: sns.scatterplot(x=error_val_pd.id, y=error_val_pd.error, hue = np.array(y_val))
plt.hlines(limit, 0, len(error_val), color='r')
```

```
Out[63]: <matplotlib.collections.LineCollection at 0x200d1d70250>
```



```
In [65]: predictions_test = autoencoder.predict(X_test)
```

```
In [66]: error_test = np.mean(np.power(X_test - predictions_test, 2), axis=1)
```