```
In [1]:   # install the opendatasets package
          !pip install opendatasets

          import opendatasets as od

          # download the dataset (this is a Kaggle dataset)
          # during download you will be required to input your Kaggle username and password
          od.download("https://www.kaggle.com/mlg-ulb/creditcardfraud")
```

Requirement already satisfied: opendatasets in c:\users\hp\anaconda3\lib\site-packages (0.1.22)
Requirement already satisfied: kaggle in c:\users\hp\anaconda3\lib\site-packages (from opendatasets) (1.5.12)
Requirement already satisfied: tqdm in c:\users\hp\anaconda3\lib\site-packages (from opendatasets) (4.62.3)
Requirement already satisfied: click in c:\users\hp\anaconda3\lib\site-packages (from opendatasets) (8.0.3)
Requirement already satisfied: colorama in c:\users\hp\anaconda3\lib\site-packages (from click->opendatasets) (0.4.4)
Requirement already satisfied: python-dateutil in c:\users\hp\anaconda3\lib\site-packages (from kaggle->opendatasets) (2.8.2)
Requirement already satisfied: requests in c:\users\hp\anaconda3\lib\site-packages (from kaggle->opendatasets) (2.26.0)
Requirement already satisfied: python-slugify in c:\users\hp\anaconda3\lib\site-packages (from kaggle->opendatasets) (5.0.2)
Requirement already satisfied: six>=1.10 in c:\users\hp\anaconda3\lib\site-packages (from kaggle->opendatasets) (1.16.0)
Requirement already satisfied: urllib3 in c:\users\hp\anaconda3\lib\site-packages (from kaggle->opendatasets) (1.26.7)
Requirement already satisfied: certifi in c:\users\hp\anaconda3\lib\site-packages (from kaggle->opendatasets) (2021.10.8)
Requirement already satisfied: text-unidecode>=1.3 in c:\users\hp\anaconda3\lib\site-packages (from python-slugify->kaggle->open
datasets) (1.3)
Requirement already satisfied: charset-normalizer~=2.0.0 in c:\users\hp\anaconda3\lib\site-packages (from requests->kaggle->open
datasets) (2.0.4)
Requirement already satisfied: idna<4,>=2.5 in c:\users\hp\anaconda3\lib\site-packages (from requests->kaggle->opendatasets) (3.
2)
Skipping, found downloaded files in ".\creditcardfraud" (use force=True to force download)

```
In [2]:   # Snap ML is available on PyPI. To install it simply run the pip command below.
          !pip install snapml
```

Collecting snapml
  Downloading snapml-1.11.1-cp39-cp39-win_amd64.whl (1.1 MB)
Requirement already satisfied: scipy in c:\users\hp\anaconda3\lib\site-packages (from snapml) (1.7.1)
Requirement already satisfied: scikit-learn in c:\users\hp\anaconda3\lib\site-packages (from snapml) (0.24.2)
Requirement already satisfied: numpy>=1.18.5 in c:\users\hp\anaconda3\lib\site-packages (from snapml) (1.20.3)
Requirement already satisfied: joblib>=0.11 in c:\users\hp\anaconda3\lib\site-packages (from scikit-learn->snapml) (1.1.0)
Requirement already satisfied: threadpoolctl>=2.0.0 in c:\users\hp\anaconda3\lib\site-packages (from scikit-learn->snapml) (2.2.
0)
Installing collected packages: snapml
Successfully installed snapml-1.11.1

```
In [3]:  # Import the libraries we need to use in this lab
         from __future__ import print_function
         import numpy as np
         import pandas as pd
         import matplotlib.pyplot as plt
         %matplotlib inline
         from sklearn.model_selection import train_test_split
         from sklearn.preprocessing import normalize, StandardScaler
         from sklearn.utils.class_weight import compute_sample_weight
         from sklearn.metrics import roc_auc_score
         import time
         import warnings
         warnings.filterwarnings('ignore')
```

```
In [4]:  # read the input data
         raw_data = pd.read_csv('creditcardfraud/creditcard.csv')
         print("There are " + str(len(raw_data)) + " observations in the credit card fraud dataset.")
         print("There are " + str(len(raw_data.columns)) + " variables in the dataset.")

         # display the first rows in the dataset
         raw_data.head()
```

```
There are 284807 observations in the credit card fraud dataset.
There are 31 variables in the dataset.
```

Out[4]:

| | Time | V1 | V2 | V3 | V4 | V5 | V6 | V7 | V8 | V9 | ... | V21 | V22 | V23 | V24 | V25 |
|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|
| 0 | 0.0 | -1.359807 | -0.072781 | 2.536347 | 1.378155 | -0.338321 | 0.462388 | 0.239599 | 0.098698 | 0.363787 | ... | -0.018307 | 0.277838 | -0.110474 | 0.066928 | 0.128539 |
| 1 | 0.0 | 1.191857 | 0.266151 | 0.166480 | 0.448154 | 0.060018 | -0.082361 | -0.078803 | 0.085102 | -0.255425 | ... | -0.225775 | -0.638672 | 0.101288 | -0.339846 | 0.167170 |
| 2 | 1.0 | -1.358354 | -1.340163 | 1.773209 | 0.379780 | -0.503198 | 1.800499 | 0.791461 | 0.247676 | -1.514654 | ... | 0.247998 | 0.771679 | 0.909412 | -0.689281 | -0.327642 |
| 3 | 1.0 | -0.966272 | -0.185226 | 1.792993 | -0.863291 | -0.010309 | 1.247203 | 0.237609 | 0.377436 | -1.387024 | ... | -0.108300 | 0.005274 | -0.190321 | -1.175575 | 0.647376 |
| 4 | 2.0 | -1.158233 | 0.877737 | 1.548718 | 0.403034 | -0.407193 | 0.095921 | 0.592941 | -0.270533 | 0.817739 | ... | -0.009431 | 0.798278 | -0.137458 | 0.141267 | -0.206010 |

5 rows × 31 columns

```
In [5]: n_replicas = 10

        # inflate the original dataset
        big_raw_data = pd.DataFrame(np.repeat(raw_data.values, n_replicas, axis=0), columns=raw_data.columns)

        print("There are " + str(len(big_raw_data)) + " observations in the inflated credit card fraud dataset.")
        print("There are " + str(len(big_raw_data.columns)) + " variables in the dataset.")

        # display first rows in the new dataset
        big_raw_data.head()
```

There are 2848070 observations in the inflated credit card fraud dataset.
There are 31 variables in the dataset.

Out[5]:

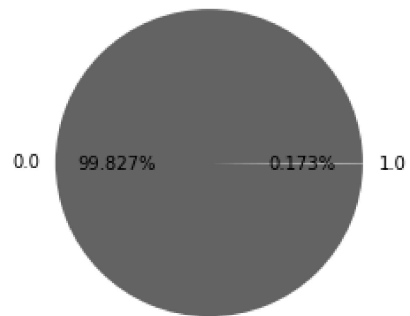| | Time | V1 | V2 | V3 | V4 | V5 | V6 | V7 | V8 | V9 | ... | V21 | V22 | V23 | V24 | V25 | |
|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|
| 0 | 0.0 | -1.359807 | -0.072781 | 2.536347 | 1.378155 | -0.338321 | 0.462388 | 0.239599 | 0.098698 | 0.363787 | ... | -0.018307 | 0.277838 | -0.110474 | 0.066928 | 0.128539 | -0.1 |
| 1 | 0.0 | -1.359807 | -0.072781 | 2.536347 | 1.378155 | -0.338321 | 0.462388 | 0.239599 | 0.098698 | 0.363787 | ... | -0.018307 | 0.277838 | -0.110474 | 0.066928 | 0.128539 | -0.1 |
| 2 | 0.0 | -1.359807 | -0.072781 | 2.536347 | 1.378155 | -0.338321 | 0.462388 | 0.239599 | 0.098698 | 0.363787 | ... | -0.018307 | 0.277838 | -0.110474 | 0.066928 | 0.128539 | -0.1 |
| 3 | 0.0 | -1.359807 | -0.072781 | 2.536347 | 1.378155 | -0.338321 | 0.462388 | 0.239599 | 0.098698 | 0.363787 | ... | -0.018307 | 0.277838 | -0.110474 | 0.066928 | 0.128539 | -0.1 |
| 4 | 0.0 | -1.359807 | -0.072781 | 2.536347 | 1.378155 | -0.338321 | 0.462388 | 0.239599 | 0.098698 | 0.363787 | ... | -0.018307 | 0.277838 | -0.110474 | 0.066928 | 0.128539 | -0.1 |

5 rows × 31 columns

```
In [6]: # get the set of distinct classes
        labels = big_raw_data.Class.unique()

        # get the count of each class
        sizes = big_raw_data.Class.value_counts().values

        # plot the class value counts
        fig, ax = plt.subplots()
        ax.pie(sizes, labels=labels, autopct='%1.3f%%')
        ax.set_title('Target Variable Value Counts')
        plt.show()
```
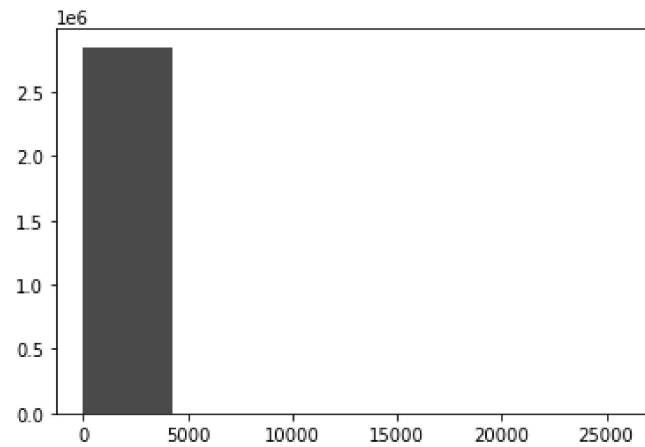
Target Variable Value Counts

0.0    99.827%    0.173%    1.0

In [7]:
```python
# we provide our solution here
plt.hist(big_raw_data.Amount.values, 6, histtype='bar', facecolor='g')
plt.show()

print("Minimum amount value is ", np.min(big_raw_data.Amount.values))
print("Maximum amount value is ", np.max(big_raw_data.Amount.values))
print("90% of the transactions have an amount less or equal than ", np.percentile(raw_data.Amount.values, 90))
```



```
Minimum amount value is  0.0
Maximum amount value is  25691.16
90% of the transactions have an amount less or equal than  203.0
```

```
In [8]:  # data preprocessing such as scaling/normalization is typically useful for
         # linear models to accelerate the training convergence

         # standardize features by removing the mean and scaling to unit variance
         big_raw_data.iloc[:, 1:30] = StandardScaler().fit_transform(big_raw_data.iloc[:, 1:30])
         data_matrix = big_raw_data.values

         # X: feature matrix (for this analysis, we exclude the Time variable from the dataset)
         X = data_matrix[:, 1:30]

         # y: labels vector
         y = data_matrix[:, 30]

         # data normalization
         X = normalize(X, norm="l1")

         # print the shape of the features matrix and the labels vector
         print('X.shape=', X.shape, 'y.shape=', y.shape)
```

```
X.shape= (2848070, 29) y.shape= (2848070,)
```

```
In [9]:  X_train, X_test, y_train, y_test = train_test_split(X, y, test_size=0.3, random_state=42, stratify=y)
         print('X_train.shape=', X_train.shape, 'Y_train.shape=', y_train.shape)
         print('X_test.shape=', X_test.shape, 'Y_test.shape=', y_test.shape)
```

```
X_train.shape= (1993649, 29) Y_train.shape= (1993649,)
X_test.shape= (854421, 29) Y_test.shape= (854421,)
```

```
In [10]: # compute the sample weights to be used as input to the train routine so that
         # it takes into account the class imbalance present in this dataset
         w_train = compute_sample_weight('balanced', y_train)

         # import the Decision Tree Classifier Model from scikit-learn
         from sklearn.tree import DecisionTreeClassifier

         # for reproducible output across multiple function calls, set random_state to a given integer value
         sklearn_dt = DecisionTreeClassifier(max_depth=4, random_state=35)

         # train a Decision Tree Classifier using scikit-learn
         t0 = time.time()
         sklearn_dt.fit(X_train, y_train, sample_weight=w_train)
         sklearn_time = time.time()-t0
         print("[Scikit-Learn] Training time (s):  {0:.5f}".format(sklearn_time))
```

```
[Scikit-Learn] Training time (s):  23.14271
```

```
In [11]:  # if not already computed,
          # compute the sample weights to be used as input to the train routine so that
          # it takes into account the class imbalance present in this dataset
          # w_train = compute_sample_weight('balanced', y_train)

          # import the Decision Tree Classifier Model from Snap ML
          from snapml import DecisionTreeClassifier

          # Snap ML offers multi-threaded CPU/GPU training of decision trees, unlike scikit-learn
          # to use the GPU, set the use_gpu parameter to True
          # snapml_dt = DecisionTreeClassifier(max_depth=4, random_state=45, use_gpu=True)

          # to set the number of CPU threads used at training time, set the n_jobs parameter
          # for reproducible output across multiple function calls, set random_state to a given integer value
          snapml_dt = DecisionTreeClassifier(max_depth=4, random_state=45, n_jobs=4)

          # train a Decision Tree Classifier model using Snap ML
          t0 = time.time()
          snapml_dt.fit(X_train, y_train, sample_weight=w_train)
          snapml_time = time.time()-t0
          print("[Snap ML] Training time (s):  {0:.5f}".format(snapml_time))
```

[Snap ML] Training time (s):  2.51329

```python
In [12]: # Snap ML vs Scikit-Learn training speedup
         training_speedup = sklearn_time/snapml_time
         print('[Decision Tree Classifier] Snap ML vs. Scikit-Learn speedup : {0:.2f}x '.format(training_speedup))

         # run inference and compute the probabilities of the test samples
         # to belong to the class of fraudulent transactions
         sklearn_pred = sklearn_dt.predict_proba(X_test)[:,1]

         # evaluate the Compute Area Under the Receiver Operating Characteristic
         # Curve (ROC-AUC) score from the predictions
         sklearn_roc_auc = roc_auc_score(y_test, sklearn_pred)
         print('[Scikit-Learn] ROC-AUC score : {0:.3f}'.format(sklearn_roc_auc))

         # run inference and compute the probabilities of the test samples
         # to belong to the class of fraudulent transactions
         snapml_pred = snapml_dt.predict_proba(X_test)[:,1]

         # evaluate the Compute Area Under the Receiver Operating Characteristic
         # Curve (ROC-AUC) score from the prediction scores
         snapml_roc_auc = roc_auc_score(y_test, snapml_pred)
         print('[Snap ML] ROC-AUC score : {0:.3f}'.format(snapml_roc_auc))
```

```
[Decision Tree Classifier] Snap ML vs. Scikit-Learn speedup : 9.21x
[Scikit-Learn] ROC-AUC score : 0.966
[Snap ML] ROC-AUC score : 0.966
```

```python
In [13]: # import the linear Support Vector Machine (SVM) model from Scikit-Learn
         from sklearn.svm import LinearSVC

         # instatiate a scikit-learn SVM model
         # to indicate the class imbalance at fit time, set class_weight='balanced'
         # for reproducible output across multiple function calls, set random_state to a given integer value
         sklearn_svm = LinearSVC(class_weight='balanced', random_state=31, loss="hinge", fit_intercept=False)

         # train a linear Support Vector Machine model using Scikit-Learn
         t0 = time.time()
         sklearn_svm.fit(X_train, y_train)
         sklearn_time = time.time() - t0
         print("[Scikit-Learn] Training time (s):  {0:.2f}".format(sklearn_time))
```

```
[Scikit-Learn] Training time (s):  54.31
```

```
In [14]:  # import the Support Vector Machine model (SVM) from Snap ML
          from snapml import SupportVectorMachine

          # in contrast to scikit-learn's LinearSVC, Snap ML offers multi-threaded CPU/GPU training of SVMs
          # to use the GPU, set the use_gpu parameter to True
          # snapml_svm = SupportVectorMachine(class_weight='balanced', random_state=25, use_gpu=True, fit_intercept=False)

          # to set the number of threads used at training time, one needs to set the n_jobs parameter
          snapml_svm = SupportVectorMachine(class_weight='balanced', random_state=25, n_jobs=4, fit_intercept=False)
          # print(snapml_svm.get_params())

          # train an SVM model using Snap ML
          t0 = time.time()
          model = snapml_svm.fit(X_train, y_train)
          snapml_time = time.time() - t0
          print("[Snap ML] Training time (s):  {0:.2f}".format(snapml_time))
```

[Snap ML] Training time (s):  7.30

```
In [15]:  # compute the Snap ML vs Scikit-Learn training speedup
          training_speedup = sklearn_time/snapml_time
          print('[Support Vector Machine] Snap ML vs. Scikit-Learn training speedup : {0:.2f}x '.format(training_speedup))

          # run inference using the Scikit-Learn model
          # get the confidence scores for the test samples
          sklearn_pred = sklearn_svm.decision_function(X_test)

          # evaluate accuracy on test set
          acc_sklearn  = roc_auc_score(y_test, sklearn_pred)
          print("[Scikit-Learn] ROC-AUC score:   {0:.3f}".format(acc_sklearn))

          # run inference using the Snap ML model
          # get the confidence scores for the test samples
          snapml_pred = snapml_svm.decision_function(X_test)

          # evaluate accuracy on test set
          acc_snapml  = roc_auc_score(y_test, snapml_pred)
          print("[Snap ML] ROC-AUC score:   {0:.3f}".format(acc_snapml))
```

[Support Vector Machine] Snap ML vs. Scikit-Learn training speedup : 7.44x
[Scikit-Learn] ROC-AUC score:   0.984
[Snap ML] ROC-AUC score:   0.985

In [16]:
```python
# get the confidence scores for the test samples
sklearn_pred = sklearn_svm.decision_function(X_test)
snapml_pred  = snapml_svm.decision_function(X_test)

# import the hinge_loss metric from scikit-learn
from sklearn.metrics import hinge_loss

# evaluate the hinge loss from the predictions
loss_snapml = hinge_loss(y_test, snapml_pred)
print("[Snap ML] Hinge loss:   {0:.3f}".format(loss_snapml))

# evaluate the hinge loss metric from the predictions
loss_sklearn = hinge_loss(y_test, sklearn_pred)
print("[Scikit-Learn] Hinge loss:   {0:.3f}".format(loss_snapml))

# the two models should give the same Hinge loss
```

[Snap ML] Hinge loss:   0.228
[Scikit-Learn] Hinge loss:   0.228

In [ ]: