

Credit Card Fraud Detection using Scikit-Learn and Snap ML

In this exercise session you will consolidate your machine learning (ML) modeling skills by using two popular classification models to recognize fraudulent credit card transactions. These models are: Decision Tree and Support Vector Machine. You will use a real dataset to train each of these models. The dataset includes information about transactions made by credit cards in September 2013 by European cardholders. You will use the trained model to assess if a credit card transaction is legitimate or not.

Import Libraries

```
# 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 /home/jupyterlab/conda/envs/python/lib/python3.7/site-packages (0.1.22)
Requirement already satisfied: tqdm in /home/jupyterlab/conda/envs/python/lib/python3.7/site-packages (from opendatasets) (4.60.0)
Requirement already satisfied: kaggle in /home/jupyterlab/conda/envs/python/lib/python3.7/site-packages (from opendatasets) (1.5.13)
Requirement already satisfied: click in /home/jupyterlab/conda/envs/python/lib/python3.7/site-packages (from opendatasets) (8.1.3)
Requirement already satisfied: importlib-metadata in /home/jupyterlab/conda/envs/python/lib/python3.7/site-packages (from click->opendatasets) (4.11.4)
Requirement already satisfied: six>=1.10 in /home/jupyterlab/conda/envs/python/lib/python3.7/site-packages (from kaggle->opendatasets) (1.16.0)
Requirement already satisfied: certifi in /home/jupyterlab/conda/envs/python/lib/python3.7/site-packages (from kaggle->opendatasets) (2022.12.7)
Requirement already satisfied: python-dateutil in /home/jupyterlab/conda/envs/python/lib/python3.7/site-packages (from kaggle->opendatasets) (2.8.2)
Requirement already satisfied: requests in /home/jupyterlab/conda/envs/python/lib/python3.7/site-packages (from kaggle->opendatasets) (2.29.0)
Requirement already satisfied: python-slugify in /home/jupyterlab/conda/envs/python/lib/python3.7/site-packages (from kaggle->opendatasets) (8.0.1)
Requirement already satisfied: urllib3 in /home/jupyterlab/conda/envs/python/lib/python3.7/site-packages (from kaggle->opendatasets) (1.26.15)
Requirement already satisfied: zipp>=0.5 in /home/jupyterlab/conda/envs/python/lib/python3.7/site-packages (from importlib-metadata->click->opendatasets) (3.15.0)
Requirement already satisfied: typing-extensions>=3.6.4 in /home/jupyterlab/conda/envs/python/lib/python3.7/site-packages (from importlib-metadata->click->opendatasets) (4.5.0)
Requirement already satisfied: text-unidecode>=1.3 in /home/jupyterlab/conda/envs/python/lib/python3.7/site-packages (from slugify->kaggle->opendatasets) (1.3)
```

```
ib/python3.7/site-packages (from python-slugify->kaggle->opendatasets) (1.3)
Requirement already satisfied: charset-normalizer<4,>=2 in /home/jupyterlab/conda/envs/python/lib/python3.7/site-packages (from requests->kaggle->opendatasets) (3.1.0)
Requirement already satisfied: idna<4,>=2.5 in /home/jupyterlab/conda/envs/python/lib/python3.7/site-packages (from requests->kaggle->opendatasets) (3.4)
Skipping, found downloaded files in "./creditcardfraud" (use force=True to force download)
```

```
In [25]: # Snap ML is available on PyPI. To install it simply run the pip command below.
!pip install snapml
```

```
Requirement already satisfied: snapml in /home/jupyterlab/conda/envs/python/lib/python3.7/site-packages (1.14.0)
Requirement already satisfied: scikit-learn in /home/jupyterlab/conda/envs/python/lib/python3.7/site-packages (from snapml) (0.20.1)
Requirement already satisfied: scipy in /home/jupyterlab/conda/envs/python/lib/python3.7/site-packages (from snapml) (1.7.3)
Requirement already satisfied: numpy>=1.18.5 in /home/jupyterlab/conda/envs/python/lib/python3.7/site-packages (from snapml) (1.21.6)
```

```
In [26]: # 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')
```

Dataset Analysis

In this section you will read the dataset in a Pandas dataframe and visualize its content. You will also look at some data statistics.

Note: A Pandas dataframe is a two-dimensional, size-mutable, potentially heterogeneous tabular data structure. For more information:

<https://pandas.pydata.org/docs/reference/api/pandas.DataFrame.html>.

```
In [27]: # 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[27]:
```

	Time	V1	V2	V3	V4	V5	V6	V7	V8
0	0.0	-1.359807	-0.072781	2.536347	1.378155	-0.338321	0.462388	0.239599	0.098698
1	0.0	1.191857	0.266151	0.166480	0.448154	0.060018	-0.082361	-0.078803	0.085102
2	1.0	-1.358354	-1.340163	1.773209	0.379780	-0.503198	1.800499	0.791461	0.247676
3	1.0	-0.966272	-0.185226	1.792993	-0.863291	-0.010309	1.247203	0.237609	0.377436
4	2.0	-1.158233	0.877737	1.548718	0.403034	-0.407193	0.095921	0.592941	-0.270533

5 rows × 31 columns

In practice, a financial institution may have access to a much larger dataset of transactions. To simulate such a case, we will inflate the original one 10 times.

```
In [28]: n_replicas = 10

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

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

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

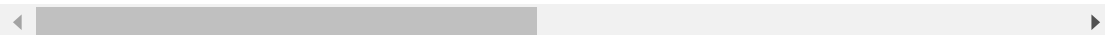
There are 2848070 observations in the inflated credit card fraud dataset.

There are 31 variables in the dataset.

```
Out[28]:
```

	Time	V1	V2	V3	V4	V5	V6	V7	V8	
0	0.0	-1.359807	-0.072781	2.536347	1.378155	-0.338321	0.462388	0.239599	0.098698	0.3637
1	0.0	-1.359807	-0.072781	2.536347	1.378155	-0.338321	0.462388	0.239599	0.098698	0.3637
2	0.0	-1.359807	-0.072781	2.536347	1.378155	-0.338321	0.462388	0.239599	0.098698	0.3637
3	0.0	-1.359807	-0.072781	2.536347	1.378155	-0.338321	0.462388	0.239599	0.098698	0.3637
4	0.0	-1.359807	-0.072781	2.536347	1.378155	-0.338321	0.462388	0.239599	0.098698	0.3637

5 rows × 31 columns



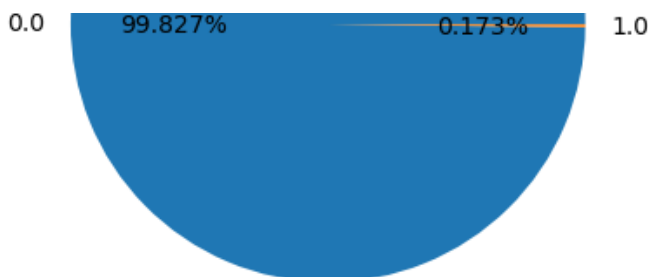
Each row in the dataset represents a credit card transaction. As shown above, each row has 31 variables. One variable (the last variable in the table above) is called Class and represents the target variable. Your objective will be to train a model that uses the other variables to predict the value of the Class variable. Let's first retrieve basic statistics about the target variable.

Note: For confidentiality reasons, the original names of most features are anonymized V1, V2 .. V28. The values of these features are the result of a PCA transformation and are numerical. The feature 'Class' is the target variable and it takes two values: 1 in case of fraud and 0 otherwise. For more information about the dataset please visit this webpage: <https://www.kaggle.com/mlg-ulb/creditcardfraud>.

```
In [29]: # get the set of distinct classes
labels = big_raw_data.Class.unique()
labels
# get the count of each class
sizes = big_raw_data.Class.value_counts().values
sizes
# 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





As shown above, the Class variable has two values: 0 (the credit card transaction is legitimate) and 1 (the credit card transaction is fraudulent). Thus, you need to model a binary classification problem. Moreover, the dataset is highly unbalanced, the target variable classes are not represented equally. This case requires special attention when training or when evaluating the quality of a model. One way of handling this case at train time is to bias the model to pay more attention to the samples in the minority class. The models under the current study will be configured to take into account the class weights of the samples at train/fit time.

Practice

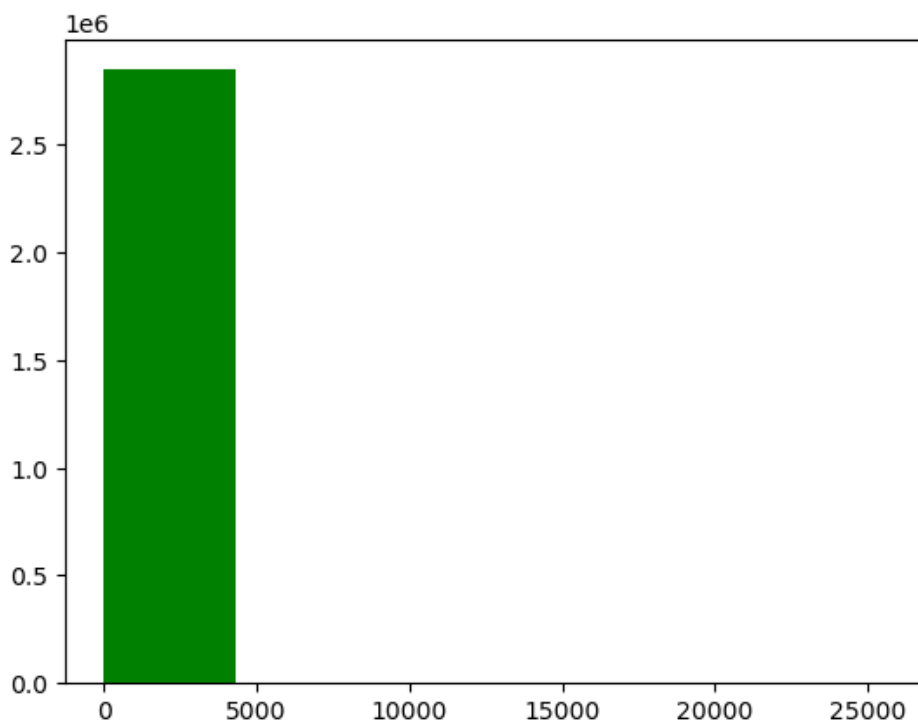
The credit card transactions have different amounts. Could you plot a histogram that shows the distribution of these amounts? What is the range of these amounts (min/max)? Could you print the 90th percentile of the amount values?

In []: `# your code here`

In [30]: `# 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_da
```



```
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
```

Dataset Preprocessing

In this subsection you will prepare the data for training.

```
In [31]: # 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,)

Dataset Train/Test Split

Now that the dataset is ready for building the classification models, you need to first divide the pre-processed dataset into a subset to be used for training the model (the train set) and a subset to be used for evaluating the quality of the model (the test set).

```
In [32]: X_train, X_test, y_train, y_test = train_test_split(X, y, test_size=0.3, random_state=42)
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,)

Build a Decision Tree Classifier model with Scikit-Learn

```
In [33]: # 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 in
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): {:.5f}".format(sklearn_time))
```

[Scikit-Learn] Training time (s): 48.20114

Build a Decision Tree Classifier model with Snap ML

```
In [34]: # 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 in
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): 9.25057

Evaluate the Scikit-Learn and Snap ML Decision Tree Classifier Models

In [19]:

```
# 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(t

# 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 : 4.64x

[Scikit-Learn] ROC-AUC score : 0.966

[Snap ML] ROC-AUC score : 0.966

As shown above both decision tree models provide the same score on the test dataset. However Snap ML runs the training routine 12x faster than Scikit-Learn. This is one of the advantages of using Snap ML: acceleration of training of classical machine learning models, such as linear and tree-based models. For more Snap ML examples, please visit [snapml-examples](#).

Build a Support Vector Machine model with Scikit-Learn

In [44]:

```
# import the Linear Support Vector Machine (SVM) model from Scikit-Learn
from sklearn.svm import LinearSVC

# instantiate 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 in
sklearn_svm = LinearSVC(class_weight='balanced', random_state=31, loss="hinge", fit_inte

# 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): 98.94

Build a Support Vector Machine model with Snap ML

In [37]:

```
# 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
# to use the GPU, set the use_gpu parameter to True
# snapml_svm = SupportVectorMachine(class_weight='balanced', random_state=25, use_gpu=True)

# 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_params={})
# print(snapml_svm.get_params())
```