# Ejercicio Random Forest

---

# Ejercicio Random Forest

---

In [1]:

```python
import pandas as pd
import numpy as np
import matplotlib.pyplot as plt
import seaborn as sns

from sklearn.metrics import confusion_matrix
from sklearn.metrics import classification_report
from sklearn.model_selection import train_test_split
from sklearn.linear_model import LogisticRegression
from sklearn.decomposition import PCA
from sklearn.tree import DecisionTreeClassifier

from pylab import rcParams

from imblearn.under_sampling import NearMiss
from imblearn.over_sampling import RandomOverSampler
from imblearn.combine import SMOTETomek
from imblearn.ensemble import BalancedBaggingClassifier

from collections import Counter

#set up graphic style in this case I am using the color scheme from xkcd.com
```

```
rcParams['figure.figsize'] = 14, 8.7 # Golden Mean
LABELS = ["Normal","Fraud"]
#col_list = ["cerulean","scarlet"]# https://xkcd.com/color/rgb/
#sns.set(style='white', font_scale=1.75, palette=sns.xkcd_palette(col_list))

%matplotlib inline
```

In [2]:
```
from google.colab import drive
drive.mount('/content/drive')
```

Mounted at /content/drive

## Cargamos Datos

In [ ]:
```
# Descargar desde https://www.kaggle.com/mlg-ulb/creditcardfraud/data

df = pd.read_csv("/content/drive/MyDrive/DATASET_TOULOUSE_C3/creditcard.csv")
df.head(n=5)
```

Out[ ]:

| | Time | V1 | V2 | V3 | V4 | V5 | V6 | V7 | V8 | V9 | ... |
|---|------|-----|-----|-----|-----|-----|-----|-----|-----|-----|-----|
| 0 | 0.0 | -1.359807 | -0.072781 | 2.536347 | 1.378155 | -0.338321 | 0.462388 | 0.239599 | 0.098698 | 0.363787 | ... | -0.01 |
| 1 | 0.0 | 1.191857 | 0.266151 | 0.166480 | 0.448154 | 0.060018 | -0.082361 | -0.078803 | 0.085102 | -0.255425 | ... | -0.22 |
| 2 | 1.0 | -1.358354 | -1.340163 | 1.773209 | 0.379780 | -0.503198 | 1.800499 | 0.791461 | 0.247676 | -1.514654 | ... | 0.24 |
| 3 | 1.0 | -0.966272 | -0.185226 | 1.792993 | -0.863291 | -0.010309 | 1.247203 | 0.237609 | 0.377436 | -1.387024 | ... | -0.10 |
| 4 | 2.0 | -1.158233 | 0.877737 | 1.548718 | 0.403034 | -0.407193 | 0.095921 | 0.592941 | -0.270533 | 0.817739 | ... | -0.00 |

5 rows × 31 columns

In [ ]:
```
df.shape
```

Out[ ]:
```
(284807, 31)
```

## Vemos Desbalanceo

In [ ]:
```
pd.value_counts(df['Class'], sort = True) #class comparison 0=Normal 1=Fraud
```

Out[ ]:
```
0    284315
1       492
Name: Class, dtype: int64
```

In [ ]:
```
normal_df = df[df.Class == 0] #registros normales
fraud_df = df[df.Class == 1] #casos de fraude
```

## Creamos Dataset

In [ ]:
```
y = df['Class']
X = df.drop('Class', axis=1)
X_train, X_test, y_train, y_test = train_test_split(X, y, train_size=0.7)
```

In [ ]:
```
def mostrar_resultados(y_test, pred_y):
    conf_matrix = confusion_matrix(y_test, pred_y)
    plt.figure(figsize=(8, 8))
    sns.heatmap(conf_matrix, xticklabels=LABELS, yticklabels=LABELS, annot=True, fmt="d"
    plt.title("Confusion matrix")
```

```
        plt.ylabel('True class')
        plt.xlabel('Predicted class')
        plt.show()
        print (classification_report(y_test, pred_y))
```

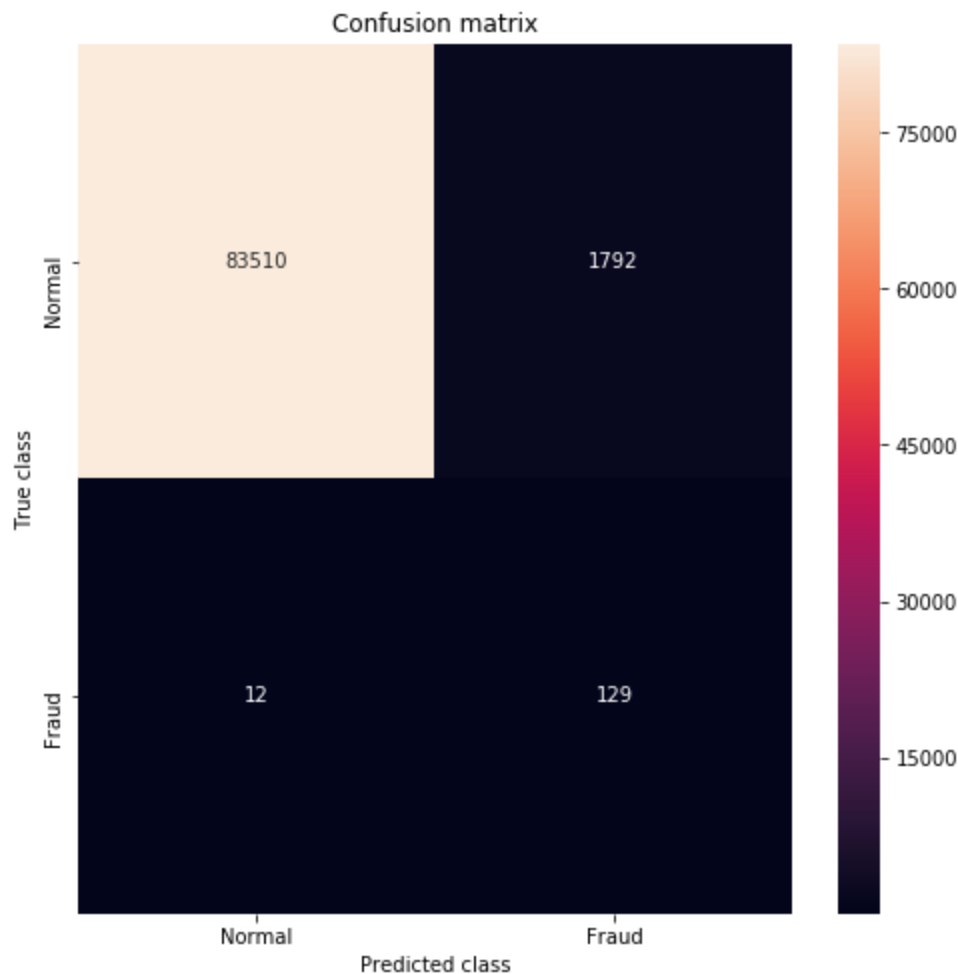## Ejecutamos Modelo con LogisticRegresion para poder Comparar

```python
In [ ]: def run_model_balanced(X_train, X_test, y_train, y_test):
            clf = LogisticRegression(C=1.0,penalty='l2',random_state=1,solver="newton-cg",class_
            clf.fit(X_train, y_train)
            return clf

        model = run_model_balanced(X_train, X_test, y_train, y_test)
```

/Users/jbagnato/anaconda3/envs/python36/lib/python3.6/site-packages/sklearn/utils/optimi
ze.py:203: ConvergenceWarning: newton-cg failed to converge. Increase the number of iter
ations.
  "number of iterations.", ConvergenceWarning)

### Veamos como responde en el test set

```python
In [ ]: pred_y = model.predict(X_test)
        mostrar_resultados(y_test, pred_y)
```



|          | precision | recall | f1-score | support |
|----------|-----------|--------|----------|---------|
| 0        | 1.00      | 0.98   | 0.99     | 85302   |
| 1        | 0.07      | 0.91   | 0.13     | 141     |
|          |           |        |          |         |
| accuracy |           |        | 0.98     | 85443   |

```
       macro avg       0.53        0.95        0.56       85443
    weighted avg       1.00        0.98        0.99       85443
```

# Probamos con Random Forest

## ATENCION: Este modelo toma algo más de tiempo en entrenar

```python
from sklearn.ensemble import RandomForestClassifier

# Crear el modelo con 100 arboles
model = RandomForestClassifier(n_estimators=100,
                               bootstrap = True,verbose=2,
                               max_features = 'sqrt')
# entrenar!
model.fit(X_train, y_train)
```

```
[Parallel(n_jobs=4)]: Using backend ThreadingBackend with 4 concurrent workers.
building tree 1 of 100building tree 2 of 100building tree 3 of 100

building tree 4 of 100

building tree 5 of 100
building tree 6 of 100
building tree 7 of 100
building tree 8 of 100
building tree 9 of 100
building tree 10 of 100
building tree 11 of 100
building tree 12 of 100
building tree 13 of 100
building tree 14 of 100
building tree 15 of 100
building tree 16 of 100
building tree 17 of 100
building tree 18 of 100
building tree 19 of 100
building tree 20 of 100
building tree 21 of 100
building tree 22 of 100
building tree 23 of 100
building tree 24 of 100
building tree 25 of 100
building tree 26 of 100
building tree 27 of 100
building tree 28 of 100
building tree 29 of 100
building tree 30 of 100
building tree 31 of 100
building tree 32 of 100
building tree 33 of 100
building tree 34 of 100
building tree 35 of 100
building tree 36 of 100

[Parallel(n_jobs=4)]: Done  33 tasks       | elapsed:   25.4s
building tree 37 of 100
building tree 38 of 100
building tree 39 of 100
building tree 40 of 100
building tree 41 of 100
```

```
        building tree 42 of 100
        building tree 43 of 100
        building tree 44 of 100
        building tree 45 of 100
        building tree 46 of 100
        building tree 47 of 100
        building tree 48 of 100
        building tree 49 of 100
        building tree 50 of 100
        building tree 51 of 100
        building tree 52 of 100
        building tree 53 of 100
        building tree 54 of 100
        building tree 55 of 100
        building tree 56 of 100
        building tree 57 of 100
        building tree 58 of 100
        building tree 59 of 100
        building tree 60 of 100
        building tree 61 of 100building tree 62 of 100

        building tree 63 of 100
        building tree 64 of 100
        building tree 65 of 100
        building tree 66 of 100
        building tree 67 of 100
        building tree 68 of 100
        building tree 69 of 100
        building tree 70 of 100
        building tree 71 of 100
        building tree 72 of 100
        building tree 73 of 100
        building tree 74 of 100
        building tree 75 of 100
        building tree 76 of 100
        building tree 77 of 100
        building tree 78 of 100
        building tree 79 of 100
        building tree 80 of 100
        building tree 81 of 100
        building tree 82 of 100
        building tree 83 of 100
        building tree 84 of 100
        building tree 85 of 100
        building tree 86 of 100
        building tree 87 of 100
        building tree 88 of 100
        building tree 89 of 100
        building tree 90 of 100
        building tree 91 of 100
        building tree 92 of 100
        building tree 93 of 100
        building tree 94 of 100
        building tree 95 of 100
        building tree 96 of 100
        building tree 97 of 100
        building tree 98 of 100
        building tree 99 of 100
        building tree 100 of 100
        [Parallel(n_jobs=4)]: Done 100 out of 100 | elapsed:  1.3min finished
```
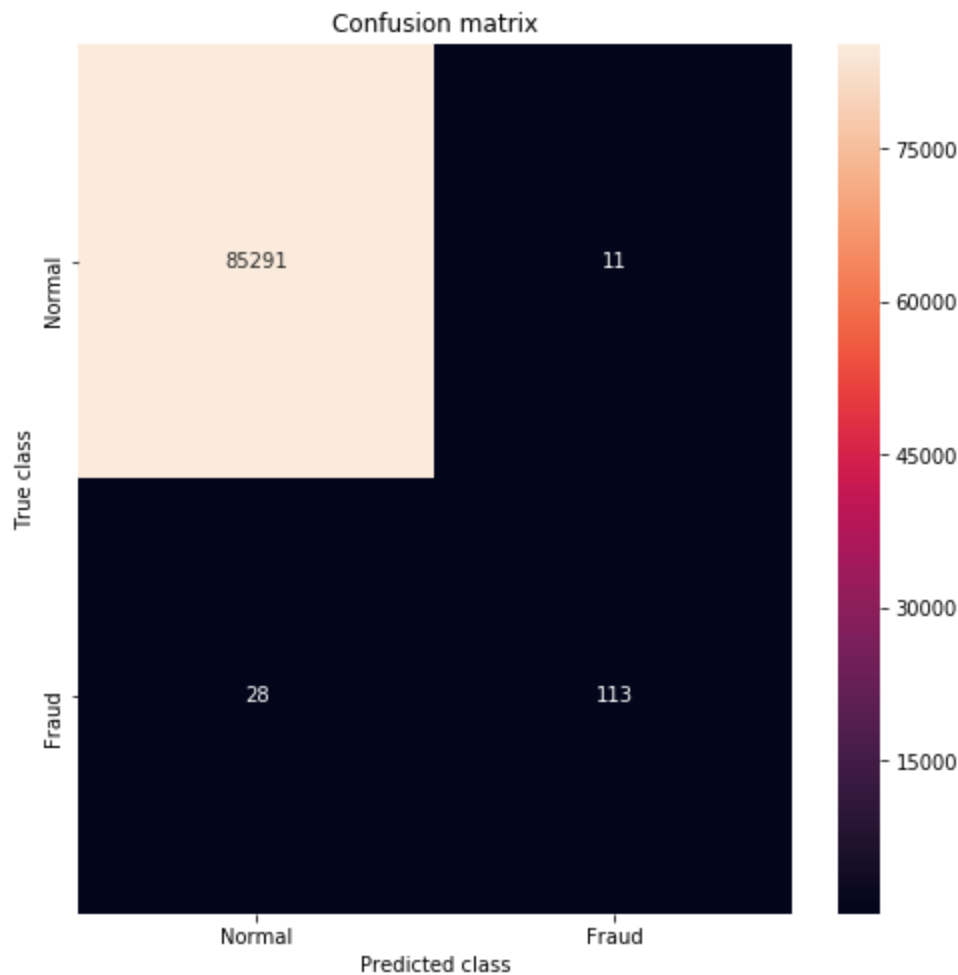
Out[ ]:
```
RandomForestClassifier(bootstrap=True, class_weight=None, criterion='gini',
                       max_depth=None, max_features='sqrt', max_leaf_nodes=None,
                       min_impurity_decrease=0.0, min_impurity_split=None,
                       min_samples_leaf=1, min_samples_split=2,
                       min_weight_fraction_leaf=0.0, n_estimators=100, n_jobs=4,
```

```
              oob_score=False, random_state=None, verbose=2,
              warm_start=False)
```

## Revisemos los resultados

```
In [ ]:  pred_y = model.predict(X_test)
         mostrar_resultados(y_test, pred_y)
```

Confusion matrix



```
              precision    recall  f1-score   support

           0       1.00      1.00      1.00     85302
           1       0.91      0.80      0.85       141

    accuracy                           1.00     85443
   macro avg       0.96      0.90      0.93     85443
weighted avg       1.00      1.00      1.00     85443
```

## Otro Bosque: Random Forest -más veloz-

```
In [ ]:  # otro modelo, variando hiperparámetros
         model = RandomForestClassifier(n_estimators=100, class_weight="balanced",
                                        max_features = 'sqrt', verbose=2, max_depth=6,
                                        oob_score=True, random_state=50)
         # a entrenar
         model.fit(X_train, y_train)
```

```
[Parallel(n_jobs=4)]: Using backend ThreadingBackend with 4 concurrent workers.
building tree 1 of 100building tree 2 of 100building tree 3 of 100
building tree 4 of 100
```

```
building tree 5 of 100
building tree 6 of 100
building tree 7 of 100
building tree 8 of 100
building tree 9 of 100
building tree 10 of 100
building tree 11 of 100
building tree 12 of 100
building tree 13 of 100
building tree 14 of 100
building tree 15 of 100
building tree 16 of 100
building tree 17 of 100
building tree 18 of 100
building tree 19 of 100
building tree 20 of 100
building tree 21 of 100
building tree 22 of 100
building tree 23 of 100
building tree 24 of 100
building tree 25 of 100
building tree 26 of 100
building tree 27 of 100
building tree 28 of 100
building tree 29 of 100
building tree 30 of 100
building tree 31 of 100
building tree 32 of 100
building tree 33 of 100
building tree 34 of 100
building tree 35 of 100
building tree 36 of 100
[Parallel(n_jobs=4)]: Done  33 tasks       | elapsed:     9.6s
building tree 37 of 100
building tree 38 of 100
building tree 39 of 100
building tree 40 of 100
building tree 41 of 100
building tree 42 of 100
building tree 43 of 100
building tree 44 of 100
building tree 45 of 100
building tree 46 of 100
building tree 47 of 100
building tree 48 of 100
building tree 49 of 100
building tree 50 of 100
building tree 51 of 100
building tree 52 of 100
building tree 53 of 100
building tree 54 of 100
building tree 55 of 100
building tree 56 of 100
building tree 57 of 100
building tree 58 of 100
building tree 59 of 100building tree 60 of 100

building tree 61 of 100
building tree 62 of 100
building tree 63 of 100
building tree 64 of 100
building tree 65 of 100
building tree 66 of 100
building tree 67 of 100
```

```
building tree 68 of 100
building tree 69 of 100
building tree 70 of 100
building tree 71 of 100
building tree 72 of 100
building tree 73 of 100
building tree 74 of 100
building tree 75 of 100
building tree 76 of 100
building tree 77 of 100
building tree 78 of 100
building tree 79 of 100
building tree 80 of 100
building tree 81 of 100
building tree 82 of 100
building tree 83 of 100
building tree 84 of 100
building tree 85 of 100
building tree 86 of 100
building tree 87 of 100
building tree 88 of 100
building tree 89 of 100
building tree 90 of 100
building tree 91 of 100
building tree 92 of 100
building tree 93 of 100
building tree 94 of 100
building tree 95 of 100
building tree 96 of 100
building tree 97 of 100
building tree 98 of 100
building tree 99 of 100
building tree 100 of 100
[Parallel(n_jobs=4)]: Done 100 out of 100 | elapsed:    28.3s finished
```
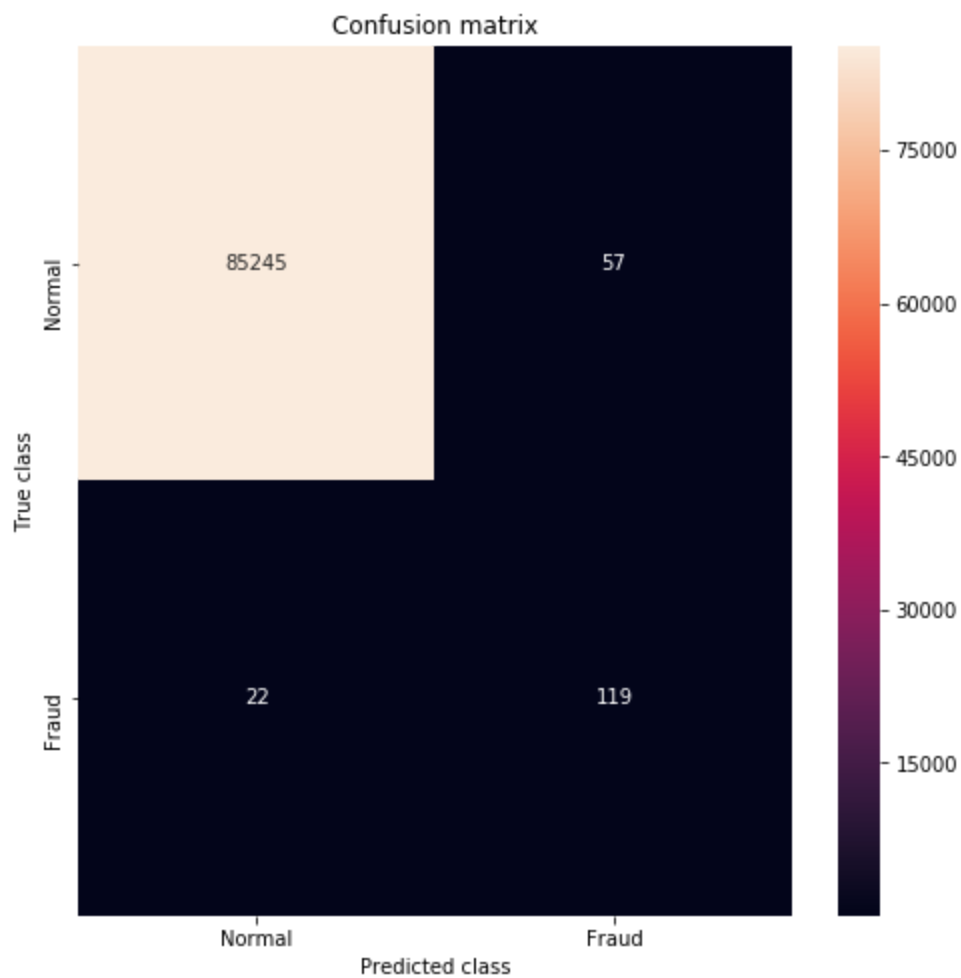
Out[ ]:
```
RandomForestClassifier(bootstrap=True, class_weight='balanced',
                       criterion='gini', max_depth=6, max_features='sqrt',
                       max_leaf_nodes=None, min_impurity_decrease=0.0,
                       min_impurity_split=None, min_samples_leaf=1,
                       min_samples_split=2, min_weight_fraction_leaf=0.0,
                       n_estimators=100, n_jobs=4, oob_score=True,
                       random_state=50, verbose=2, warm_start=False)
```

## Veamos la Confusion Matrix con el conjunto de Test

In [ ]:
```python
pred_y = model.predict(X_test)
mostrar_resultados(y_test, pred_y)
```

```
[Parallel(n_jobs=4)]: Using backend ThreadingBackend with 4 concurrent workers.
[Parallel(n_jobs=4)]: Done  33 tasks      | elapsed:    0.1s
[Parallel(n_jobs=4)]: Done 100 out of 100 | elapsed:    0.3s finished
```

Confusion matrix

```
              precision    recall  f1-score   support

           0       1.00      1.00      1.00     85302
           1       0.68      0.84      0.75       141

    accuracy                           1.00     85443
   macro avg       0.84      0.92      0.88     85443
weighted avg       1.00      1.00      1.00     85443
```

# Comprobamos Resultados

```python
In [ ]:  from sklearn.metrics import roc_auc_score

         # Calculate roc auc
         roc_value = roc_auc_score(y_test, pred_y)
```

```python
In [ ]:  print(roc_value)

         0.9216517085479026
```

El valor de roc cuanto más cerca de 1, mejor. si fuera 0.5 daría igual que fuesen valores aleatorios y sería un mal modelo