

DATOS MASIVOS

Shamir Alejandro Rodriguez Ojeda

Matricula: 1423061

Herramienta: Colab

```

1 #####
2 ###                               PIA - SHAMIR A. RODRIGUEZ OJEDA          ##
3 ###                               Maestria en Ciencia de Datos           ##
4 ###                               12/04/2024                             ##
5 ###                               ##
6 #####

```

✓ Instalación Librerías

```

[ ] 1 !pip install --q https://github.com/pandas-profiling/pandas-profiling/archive/master.zip
2 !pip install --q pyod # Detección de atípicos +30 algoritmos
3 !pip install --q catboost # Catboost
4 !pip install --q imblearn # Balanceo de clases

```

```

- 17.8 MB 26.6 MB/s 0:00:01
Preparing metadata (setup.py) ... done
104.8/104.8 kB 1.4 MB/s eta 0:00:00
Preparing metadata (setup.py) ... done
686.1/686.1 kB 7.5 MB/s eta 0:00:00
296.5/296.5 kB 5.4 MB/s eta 0:00:00
Building wheel for ydata-profiling (setup.py) ... done
Building wheel for htmlmin (setup.py) ... done
165.0/165.0 kB 4.1 MB/s eta 0:00:00
Preparing metadata (setup.py) ... done
Building wheel for pyod (setup.py) ... done
98.2/98.2 MB 8.7 MB/s eta 0:00:00

```

```

[ ] 1 from google.colab import drive
2 from pathlib import Path

```

```

[ ] 1 import numpy as np
2 import pandas as pd
3 from pandas_profiling import ProfileReport

```

```
[ ] 1 from google.colab import drive
2 from pathlib import Path

[ ] 1 import numpy as np
2 import pandas as pd
3 from pandas_profiling import ProfileReport

⚠️ <ipython-input-3-c83a24b5f2b>:3: DeprecationWarning: 'import pandas_profiling' is going to be deprecated by April 1st. Please use 'import ydata_profiling' instead.
from pandas_profiling import ProfileReport

[ ] 1 #Drive.mount('/content/drive', force_reconnect=True)

⚠️ Mounted at /content/drive

1 from pyod.models.lof import LOF
2 from sklearn.utils.class_weight import compute_class_weight
3 from sklearn.metrics import classification_report
4 from sklearn.metrics import roc_auc_score
5 from sklearn.metrics import confusion_matrix
6 from sklearn.metrics import average_precision_score
7 from sklearn import preprocessing
8 from sklearn.model_selection import GridSearchCV
9
10
11 from sklearn.model_selection import train_test_split
12 from catboost import CatBoostClassifier
13 from imblearn.over_sampling import ADASYN
14 from collections import Counter
15
16 from sklearn.model_selection import StratifiedKFold
17 from sklearn.model_selection import cross_val_score
18 from catboost import CatBoostClassifier, Pool
19
20
21 import seaborn as sns
22 import matplotlib.pyplot as plt
23 from string import ascii_letters
24
```

```
[ ] 20
21 import seaborn as sns
22 import matplotlib.pyplot as plt
23 from string import ascii_letters
24
25 import warnings
26 warnings.filterwarnings("ignore")
```

```
[ ] 1 SEED=111
```

✓ Funciones

```
[ ] 1 def grafica_correla(df, paleta="Greens"):
2     import seaborn as sns
3     import matplotlib.pyplot as plt
4     from string import ascii_letters
5     fig, ax = plt.subplots(figsize=(6, 6))
6     mask = np.zeros_like(df.corr())
7     mask[np.triu_indices_from(mask)] = 1
8     sns.heatmap(df.corr(), mask= mask, cmap=paleta)
```

```
[ ] 1 def refusion_dimension(df, umbral=0.95):
2     corr_df = df.corr().abs()
3     mask = np.triu(np.ones_like(corr_df, dtype=bool))
4     correlacion = df.corr()
5     tri_df = correlacion.mask(mask)
6     to_drop = [c for c in tri_df.columns if any(tri_df[c] > umbral)]
7     return to_drop
```

Para poder realizar la entrega siguiendo el formato del archivo TRAIN, se requiere saber las variables iniciales dadas para el caso de uso.

```
[ ] 1 etiquetas= train.iloc[:,1-7].columns
2 len(etiquetas)
```

27

Se creo un reporte previo para la exploración de los datos EDA. El archivo html tiene varios analisis sobre los datos iniciales.

```
[ ] 1 profile = ProfileReport(train, title='Pandas Profiling Report', html={'style':{'full_width':False}})
```

```
[ ] 1 profile.to_notebook_iframe()
```

```
[ ] 1 profile.to_file(output_file=link_lectura+name_report)
```

Export report to file: 100%  1/1 [00:00<00:00, 1.36Mb]

1 profile

Summarize dataset: 100%  619/619 [02:47<00:00, 2.18Mb, Completed]

Generate report structure: 100%  1/1 [00:23<00:00, 23.02s/1]

Render HTML: 100%  1/1 [00:29<00:00, 29.99s/1]

Vale la pena verificar cuantos nulos tiene el dataframe dado que muchas veces es/será un dolor de cabeza para la creación de un modelo de ML

```
[ ] 1 train.isnull().sum()
```

```
X_Minimum      0
X_Maximum      0
Y_Minimum      0
Y_Maximum      0
Pixels_Areas    0
X_Perimeter     0
Y_Perimeter     0
Sum_of_Luminosity 0
Minimum_of_Luminosity 0
Maximum_of_Luminosity 0
Length_of_Conveyer 0
TypeOfSteel_A300 0
TypeOfSteel_A400 0
Steel_Plate_Thickness 0
Edges_Index     0
Empty_Index     0
Square_Index    0
Outside_X_Index 0
Edges_X_Index   0
Edges_Y_Index   0
Outside_Global_Index 0
LogOfAreas      0
Log_X_Index     0
Log_Y_Index     0
Orientation_Index 0
Luminosity_Index 0
SigmoidOfAreas  0
Pastry          0
Z_Scratch       0
K_Scratch       0
Stains          0
Dirtiness       0
Bumps           0
Other_Faults    0
dtype: int64
```

```
1. train.describe()
```

	X_Minimum	X_Maximum	Y_Minimum	Y_Maximum	Pixels_Areas	X_Perimeter	Y_Perimeter	Sum_of_Luminosity	Minimum_of_Luminosity	Maximum_of_Luminosity	...	Orientatoin_Index	Luminosity_Index	SigmoidHArea	Po
count	1455.000000	1455.000000	1.455000e+03	1.455000e+03	1455.000000	1455.000000	1455.000000	1.455000e+03	1455.000000	1455.000000	...	1455.000000	1455.000000	1455.000000	1455.00
mean	675.056352	676.383188	1.675284e+08	1.676050e+08	1568.096252	115.644674	85.914389	2.589553e+08	84.561198	130.021993	...	0.087681	-0.130914	0.562273	0.08
std	516.641171	483.543339	1.806262e+06	1.806270e+06	5545.548258	331.232886	467.966747	8.333202e+05	32.375884	18.461963	...	0.499667	0.151786	0.340529	-0.27
min	0.000000	4.000000	8.712000e+03	8.724000e+03	2.000000	2.000000	1.000000	2.500000e+02	0.000000	37.000000	...	-0.973300	-0.968500	0.119000	0.00
25%	36.000000	194.000000	5.020949e+05	5.020949e+05	82.500000	15.000000	13.000000	9.374000e+03	83.000000	124.000000	...	-0.292200	-0.198200	0.245300	-0.00
50%	439.000000	466.000000	1.223950e+08	1.223950e+08	172.000000	26.000000	25.000000	1.888500e+04	90.000000	127.000000	...	0.098900	-0.132600	0.496300	0.00
75%	1646.500000	1075.000000	2.183653e+06	2.183653e+06	794.000000	82.000000	83.000000	8.188950e+04	106.000000	140.000000	...	0.809500	-0.086100	0.988800	0.00
max	1785.000000	1713.000000	1.296766e+07	1.296766e+07	152655.000000	10449.000000	16152.000000	1.158145e+07	203.000000	247.000000	...	0.991700	0.642100	1.000000	1.00

8 rows x 16 columns

```
1. train.head()
```

	X_Minimum	X_Maximum	Y_Minimum	Y_Maximum	Pixels_Areas	X_Perimeter	Y_Perimeter	Sum_of_Luminosity	Minimum_of_Luminosity	Maximum_of_Luminosity	...	Orientatoin_Index	Luminosity_Index	SigmoidHArea	Fasto	Scratch
1309	1343	1351	46714	46725	36	11	11	3729	30	117	...	0.2723	-0.2007	0.1934	0	9
1132	354	370	409908	409918	116	19	19	12174	82	126	...	-0.3750	-0.1801	0.2732	0	9
1230	1056	1077	4364066	4364076	106	35	32	10747	89	116	...	-0.0236	-0.2679	0.5009	0	9
783	816	822	899658	899663	18	8	6	2364	127	141	...	0.0006	-0.0347	0.1449	6	9
1762	78	94	1523381	1523431	463	69	54	60018	118	141	...	0.6900	-0.0212	0.9690	0	9

5 rows x 16 columns

Se va a proceder a eliminar las clases individuales del DF train para despues concatenar la columna target

```
1. train.drop(mask, axis=1, inplace=True)
2. train
```

	X_Minimum	X_Maximum	Y_Minimum	Y_Maximum	Pixels_Areas	X_Perimeter	Y_Perimeter	Sum_of_Luminosity	Minimum_of_Luminosity	Maximum_of_Luminosity	...	Outside_X_Index	Edges_X_Index	Edges_Y_Index
1309	1343	1351	46714	46725	36	11	11	3729	30	117	...	0.0059	0.7273	1.0000
1132	354	370	409908	409918	116	19	19	12174	82	126	...	0.0118	0.8421	1.0000
1230	1056	1077	4364066	4364076	106	35	32	10747	89	116	...	0.0154	0.6000	0.4545
783	816	822	899658	899663	18	8	6	2364	127	141	...	0.0044	0.7500	1.0000
1762	78	94	1523381	1523431	463	69	54	60018	118	141	...	0.0116	0.2319	0.9259
...
555	677	686	3379334	3379344	69	12	10	7774	95	132	...	0.0054	0.7500	1.0000
1702	114	134	115356	115418	532	75	73	54451	94	111	...	0.0144	0.2667	0.6493
194	109	121	560425	560445	124	18	20	19549	97	124	...	0.0089	0.9667	1.0000
1305	573	366	1058240	1058279	308	29	39	37264	113	132	...	0.0095	0.4483	1.0000
527	39	178	2627229	2627391	7998	556	343	889727	49	140	...	0.0993	0.2491	0.4723

1450 rows x 15 columns

Reversa de One-Hot Encoding

Se crea la columna target que toma el one-hot encoding y se hace el proceso de reversa.

```
[ ] 1 train['target'] = classes.apply(lambda x: x.idxmax(), axis = 1)
2 train
```

	A_Minimum	X_Maximum	V_Minimum	V_Maximum	Plots_Areas	X_Perimeter	V_Perimeter	Sum_of_Laminicity	Minimum_of_Laminicity	Maximum_of_Laminicity	...	Edges_X_Index	Edges_V_Index	Outside_Global_Index	logOfArea	log_X_Index
1309	1343	1351	46714	46725	56	11	11	5729	90	117	...	0.7273	1.0000	1.0	1.7482	0.9031
1132	354	370	409908	409918	116	19	10	12174	82	126	...	0.8421	1.0001	0.0	2.0645	1.3041
1336	1006	1077	4364956	4364076	106	35	22	10747	89	116	...	0.6000	0.4545	0.0	2.0253	1.3222
750	816	822	809629	809665	18	8	8	2384	127	141	...	0.7500	1.0000	0.5	1.2553	0.7782
1762	78	54	1523391	1523431	463	69	54	60618	118	141	...	0.2319	0.5059	1.0	2.0656	1.2041
...
955	677	686	5379334	5379344	89	12	10	7774	95	132	...	0.7500	1.0000	5.0	1.6388	0.9342
1792	114	134	115396	115418	582	70	78	54451	94	111	...	0.2667	0.8493	1.0	2.7255	1.3010
194	109	121	560425	560445	128	18	20	13549	97	128	...	0.8667	1.0001	1.0	2.0334	1.0792
1335	373	386	1058240	1058279	308	29	29	37264	113	132	...	0.4483	1.0000	1.0	2.4085	1.1139
627	39	178	2527229	2527391	7988	108	343	889727	49	140	...	0.2491	0.4723	1.0	3.9030	2.1430

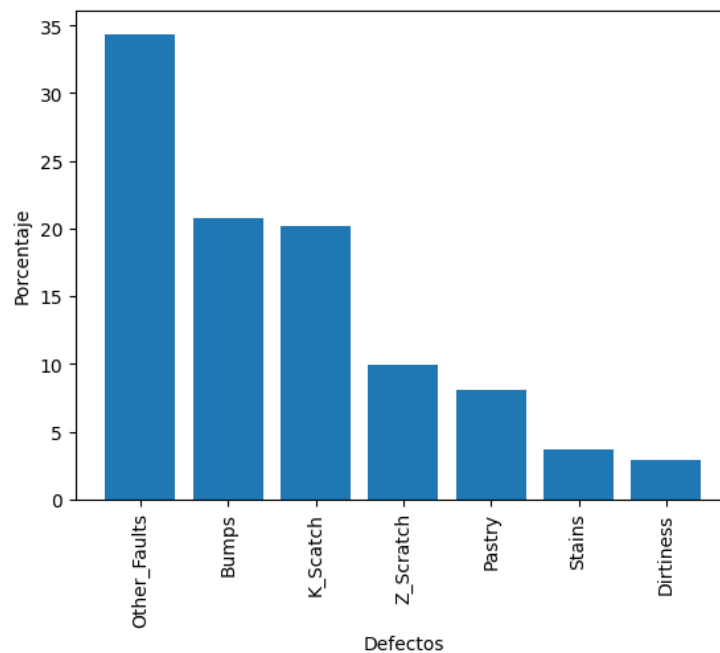
1455 rows x 26 columns

```
[ ] 1 aux
```

	class	count
0	Other_Faults	34364261
1	Bumps	20756014
2	K_Scratch	20206186
3	Z_Scratch	9968636
4	Pastry	8109966
5	Stains	3711340
6	Dirtyess	2886998

```
1 aux = ((train['target'].value_counts()/train.shape[0])*100).rename_axis('class').reset_index()
2 plt.bar(aux['class'], aux['count'],)
3 plt.xticks(rotation=90)
4 plt.ylabel('Porcentaje')
5 plt.xlabel('Defectos')
```

Text(0.5, 0, 'Defectos')



▼ Nulos

En celdas anteriores como previo a la carga de los datos se exploró rápidamente los valores nulos ahora se procederá a revisar si existe al menos 1 en todo el DF.

```
[ ] 1 train.isnull().sum().sum()
```

```
0
```

Uno de los problemas típicos al trabajar con datos es que los datos vengan en otro formato, es decir cuando uno ve el dato es un número 1 pero el sistema lo ha leído como texto '1' causando complicaciones en los modelos.

▼ formato de las variables (tipo)

```
[ ] 1 train.dtypes
```

```
X_Minimum      int64
X_Maximum      int64
V_Minimum      int64
V_Maximum      int64
Pixels_Areas    int64
K_Perimeter     int64
V_Perimeter     int64
Sum_of_Luminosity int64
Minimum_of_Luminosity int64
Maximum_of_Luminosity int64
Length_of_Conveyer int64
TypeOfSteel_A388 int64
TypeOfSteel_A488 int64
Steel_Plate_Thickness int64
Edges_Index     float64
Empty_Index     float64
Square_Index    float64
Outside_X_Index float64
Edges_X_Index   float64
```

▼ Duplicados

Validación de los duplicados en el DF

```
[ ] 1 aux=train.shape[0]
2 train.drop_duplicates(keep='first',inplace=True)
3 aux=train.shape[0]
```

```
True
```

Al parecer no existe duplicados en el DF, dado que la cantidad de columnas antes y después de hacer la operación sigue siendo la misma

▼ Atípicos

Local Outlier Factor

```
[ ] 1 clf = LOF(contamination=0.1)
2 clf.fit(train.drop('target',axis=1))
3 atipicos = clf.predict(train.drop('target',axis=1)) ## etiquetas binarias (0: no atípicos, 1: atípicos)
4
5 atipicos[atipicos!=1]
6 train=train.iloc[atipicos,:]
7 train.shape
```

```
(1322, 28)
```

▼ Alta Correlación

```
[ ] 1 correlacion = train.drop('target',axis=1).corr().unstack()
2 correla_ordena= correlacion.sort_values(kind="quicksort")
3
```

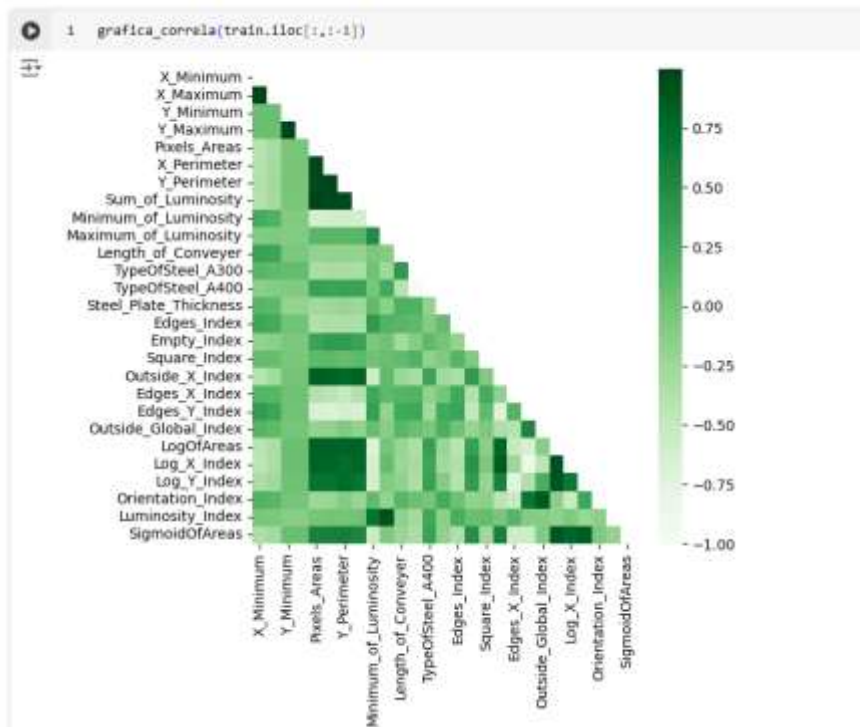

↘ Alta Correlación

```

1 correlacion = train.drop('target',axis=1).corr().unstack()
2 correla_ordena= correlacion.sort_values(kind="quicksort")
3
4 correla_ordena=correla_ordena[(correla_ordena > 0.90) & ( correla_ordena<0.999)]
5 correla_ordena
  
```

```

LogOfAreas      Log_X_Index      0.901216
Log_X_Index      LogOfAreas      0.901216
Maximum_of_Luminosity  Luminosity_Index  0.913880
Luminosity_Index  Maximum_of_Luminosity  0.913880
Y_Perimeter      Pixels_Areas      0.954729
Pixels_Areas      Y_Perimeter      0.954729
Y_Perimeter      Sum_of_Luminosity  0.958623
Sum_of_Luminosity  Y_Perimeter      0.958623
X_Perimeter      Pixels_Areas      0.974513
Pixels_Areas      X_Perimeter      0.974513
Sum_of_Luminosity  X_Perimeter      0.976020
X_Perimeter      Sum_of_Luminosity  0.976020
Y_Perimeter      X_Perimeter      0.977837
X_Perimeter      Y_Perimeter      0.977837
X_Maximum        X_Minimum        0.991427
X_Minimum        X_Maximum        0.991427
Sum_of_Luminosity  Pixels_Areas      0.997331
Pixels_Areas      Sum_of_Luminosity  0.997331
dtype: float64
  
```




```
[ ] 1 to_drop = refusion_dimension(train.drop('target',axis=1))
    2 to_drop
```

```
⇒ ['X_Minimum', 'Y_Minimum', 'Pixels_Areas', 'X_Perimeter', 'Y_Perimeter']
```

```
[ ] 1 train = train.drop(to_drop, axis=1)
    2 train.shape
```

```
⇒ (1322, 23)
```

```
[ ] 1 X=train.drop('target', axis=1).values
    2 y=train['target'].values
    3
    4 X.shape, y.shape
```

```
⇒ ((1322, 22), (1322,))
```

▼ Transformación de datos

Para poder trabajar con algunos modelos de ML se quiere hacer un encoding para las variables, normalización/estandarización segun corresponda.

```
[ ] 1 label_encoder = preprocessing.LabelEncoder()
    2 y=label_encoder.fit_transform(y)
    3 y
```

```
⇒ array([3, 0, 0, ..., 6, 3, 2])
```

```
[ ] 1 X=preprocessing.scale(X)
    2 X
```

```
⇒ array([[ 1.45284097, -0.90049171, -0.41018559, ...,  0.4166186 ,
          -0.46622624, -1.05573526],
        [-0.52164727, -0.0998491 , -0.39501208, ..., -0.92790963,
          -0.3318459 , -0.79539686],
        [ 0.90135292,  1.48458806, -0.39037167, ..., -1.23679615,
          -0.51319413, -0.57784136],
        ...,
        [-1.02281707, -0.61669185, -0.3917749 , ...,  0.68087434,
          -0.1120101 , -0.4358386 ],
        [-0.48944359, -0.34166807, -0.33594248, ...,  1.23450361,
          0.48552579,  0.77452253],
        [-0.90809145,  0.52517332,  1.67101813, ...,  0.14530407,
          -0.01089868,  1.39168837]])
```

```
[ ] 1 X_train, X_test, y_train, y_test = train_test_split(
    2     X, y, test_size=0.2, random_state=SEED, stratify=y, shuffle=True)
```

▼ Catboost

Se usará una estrategia para que Catboost sea la proporción que tienen las clases

```
[ ] 1 classes = np.unique(y_train)
    2 weights = compute_class_weight(class_weight='balanced', classes=classes, y=y_train)
    3 class_weights = dict(zip(classes, weights))
```

```
[ ] 1 train_dataset = Pool(X_train, y_train)
    2
    3 rscv = StratifiedKFold(n_splits=3, shuffle=True, random_state=SEED)
    4
    5 model = CatBoostClassifier(
    6     learning_rate=0.01
    7     ,depth=6
    8     ,n_estimators=100
    9     ,loss_function='MultiClass'
10     ,train_dir='crossentropy'
11     ,early_stopping_rounds=200
12     ,bagging_temperature = 1
13     ,metric_period = 100
14     ,langevin=True
15     ,random_seed =SEED
16     ,class_weights =class_weights
17
18 )
19
20 model.fit(train_dataset, verbose=False, plot=True)
```

<catboost.core.CatBoostClassifier at 0x79932acf3880>

▼ Datos de Entrenamiento

```
[ ] 1 y_pred = model.predict(X_train)
    2 print(classification_report(y_train,y_pred))
```

```

precision    recall  f1-score   support

0           0.62      0.71      0.66       235
1           0.39      1.00      0.56        28
2           0.95      0.89      0.92       199
3           0.77      0.37      0.50       362
4           0.39      0.80      0.53        85
5           0.76      0.91      0.83        43
6           0.71      0.90      0.79       105

accuracy          0.67      1057
macro avg          0.66      0.80      0.68      1057
weighted avg       0.72      0.67      0.66      1057
```

▼ Datos de test

```
[ ] 1 y_pred = model.predict(X_test)
    2 print(classification_report(y_test,y_pred))
```

```

precision    recall  f1-score   support

0           0.55      0.71      0.62        59
1           0.35      1.00      0.52         7
2           1.00      0.82      0.90        50
3           0.75      0.30      0.43        91
4           0.36      0.62      0.46        21
5           0.73      1.00      0.85        11
6           0.60      0.92      0.73        26

accuracy          0.62      265
macro avg          0.62      0.77      0.64      265
weighted avg       0.69      0.62      0.61      265
```

▼ HyperTunning para Catboost

▼ Balanceo de clases

Se usara una libreria para hacer el balanceo de clases, se probó una estrategia de decirle al algoritmo la proporción de las clases y se esperaba que este resultado fuera bueno. Dado que el supuesto es falso se procedió a hacer un balanceo de clases con estrategia más robustas.

```
1 oversample = ADASYN(n_neighbors=6, random_state=SEED)
2 X_balance, y_balance = oversample.fit_resample(X_test, y_test)
3
4 clases_balance= pd.DataFrame.from_dict(Counter(y_balance), orient='index').reset_index()
5 clases_balance.columns=['defectos', 'porcentaje']
6 clases_balance['porcentaje']=round((clases_balance['porcentaje']/len(X_balance))*100,2)
7 clases_balance
```

	defectos	porcentaje
0	2	14.37
1	3	14.22
2	0	15.62
3	4	14.37
4	5	13.12
5	1	14.06
6	5	14.22

Cuando se entreno el modelo baseline usando CatBoost se uso una estrategia para decirle al algoritmo como estaban distribuidas las clases. Para hacer el hypertunning no vamos a usar está estrategia sino que por el contrario usaremos una estrategia de balanceo de clases por lo cual es conveniente que el modelo base no tenga está opción impuesta.

```
1 model = CatBoostClassifier(iterations=500,
2                             learning_rate=0.01
3                             ,depth=6
4                             ,n_estimators=100
5                             ,loss_function='MultiClass'
6                             ,train_dir= 'crossentropy'
7                             ,early_stopping_rounds=200
8                             ,metric_period = 100
9                             ,random_seed =SEED
10                            ,task_type="GPU"
11                            )
12
13 model.fit(train_dataset, verbose=False, plot=True)
14
15 parameters = {'depth' : [5,6,7,9, 10],
16               'learning_rate' : [0.01,0.02,0.03,0.04],
17               'n_estimators' : [75,100,120]
18               }
19
20 Grid_CBC_balance = GridSearchCV(estimator=model
21                                 ,param_grid = parameters
22                                 ,cv = rscv
23                                 ,scoring="accuracy"
24                                 ,n_jobs=-1)
25
26 Grid_CBC_balance.fit(X_balance, y_balance)
```

```
0:   learn: 1.8742584      total: 48.8ms  remaining: 5.8s
100: learn: 0.4070332    total: 6.79s  remaining: 1.28s
119: learn: 0.3480472    total: 7.59s  remaining: 0us
```

```
> GridSearchCV
> estimator: CatBoostClassifier
> CatBoostClassifier
```

```
[ ] 1 from google.colab import output
    2 output.enable_custom_widget_manager()
```

Support for third party widgets will remain active for the duration of the session. To disable support:

```
[ ] 1 from google.colab import output
    2 output.disable_custom_widget_manager()
```

▼ Datos de Test

```
[ ] 1 Grid_CBC_balance.best_params_,Grid_CBC_balance.best_estimator_
```

```
{'depth': 10, 'learning_rate': 0.04, 'n_estimators': 120},
<catboost.core.CatBoostClassifier at 0x7f2af28cd650>
```

```
[ ] 1 y_pred = Grid_CBC_balance.best_estimator_.predict(X_test)
    2 print(classification_report(y_test,y_pred))
```

```
precision    recall  f1-score   support

0           0.96         0.99         0.97         88
1           0.85         1.00         0.92         11
2           1.00         1.00         1.00         75
3           1.00         0.94         0.97        136
4           0.94         1.00         0.97         32
5           1.00         1.00         1.00         16
6           0.97         1.00         0.99         39

accuracy          0.98         0.98         0.98        397
macro avg          0.96         0.99         0.97        397
weighted avg          0.98         0.98         0.98        397
```

▼ Validación del modelo

Ahora que se tiene el modelo entrenado, se usarán el otro conjunto de datos suministrados para validarlo.

```
[ ] 1 test=pd.read_csv(link_lectura+name_file_test, sep='\t',index_col=0)
    2 print("La dimension del DataFrame de test es {}".format(test.shape))
```

```
-----
NameError                                Traceback (most recent call last)
<ipython-input-1-c5b95c631370> in <module>()
----> 1 test=pd.read_csv('/content/drive/My Drive/Colab Notebooks/Fecnoap/Faults.test', sep='\t',index_col=0)
      2 print("La dimension del DataFrame de test es {}".format(test.shape))

NameError: name 'pd' is not defined
```

```
[ ] 1 test
```

```
-----
NameError                                Traceback (most recent call last)
<ipython-input-2-4e1243b622c6> in <module>()
----> 1 test

NameError: name 'test' is not defined
```

Se coloca un procedimiento para garantizar que el nombre de las columnas en los dos archivos entregados por la empresa sean similares, excluyendo de DF de train la variable objetivo.

```
[ ] 1 test= test[etiquetas]
    2 test
```

```
[ ] 1 XX_test=test.values
    2 XX_test.shape
```

```
[ ] 1 XX_test=preprocessing.scale(XX_test)
```

```
▶ 1 yy_pred = Grid_CBC_balance.best_estimator_.predict(XX_test)
  2 defectos= label_encoder.inverse_transform(yy_pred)
  3 defectos
```

```
[ ] 1 test['target']=defectos
    2 test
```

```
[ ] 1 test.to_csv(link_lectura+name_user+"/"+name_file_exist)
```