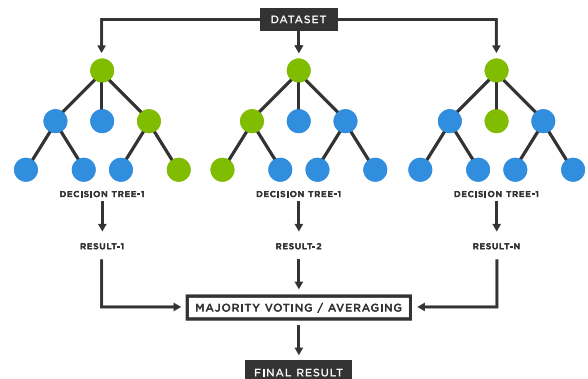


Modelado y Evaluación

Técnica del modelado

Para nuestro proyecto, era muy importante que de todos los contratos que ha adjudicado el estado se clasificara que contratos fueron cerrados y de esta forma poder realizar un modelo que pudiera predecir ésta clasificación **cerrado** según los parámetros definidos para su entrenamiento.

En Machine Learning se cuentan con varios algoritmos de clasificación de aprendizaje supervisado, tales como K-Nearest Neighbors, Regresión Logística, Naive Bayes, Árboles de decisiones, Random Forest y Redes Neuronales. Para nuestro proyecto decidimos utilizar Random Forest como algoritmo de clasificación



```
from sklearn.model_selection import train_test_split
```

debido a que nos ofrece la clasificación árboles de decisión que participan en si para la clasificación de un determinado elemento.

```
train_x, rest_x, train_y, rest_y = train_test_split(contratos_data, estado_contrato, train_size=training_count)
test_x, validate_x, test_y, validate_y = train_test_split(rest_x, rest_y, train_size=test_count)
```

Conjuntos de Entrenamiento, Pruebas y Validación

Una vez tratado el archivo original de los contratos del SECOP II, obtuvimos un DataFrame de 549.407 registros con 34 parámetros; de este DataFrame decidimos tomar el 60% de los registros para entrenamiento, el 20% para pruebas y el otro 20% para validación.

```
# Calculate test and validation set size:
original_count = len(df_final)
training_size = 0.60 # 60% of records
test_size = (1 - training_size) / 2

training_count = int(original_count * training_size)
test_count = int(original_count * test_size)
validation_count = original
```

Conjunto de datos y sus registros

Conjunto de entrenamiento	329.644
Conjunto de pruebas	109.881
Conjunto de Validación	109.882

```
_count - training_count - t  
est_count
```

Evaluación del Modelo

Para la evaluación del modelo utilizamos dos métricas de clasificación para medir que tan bien hace las predicciones con los parámetros definidos para el nuestro algoritmo de clasificación Random Forest.

Accuracy

Usamos esta métrica para determinar la exactitud de las predicciones correctas sobre el total de predicciones realizadas

Recall

Utilizamos esta métrica para medir la capacidad del modelo de identificar todos los casos positivos (verdaderos positivos)

Análisis de entrenamiento del modelo

Inicialmente validamos nuestro modelo usando 100 estimadores (# de árboles) y tomando los demás parámetros por default del algoritmo

```
from sklearn.ensemble import RandomForestClassifier  
from sklearn.svm import LinearSVC  
  
model = RandomForestClassifier(n_estimators=100)  
model.fit(features_train_x, train_y)  
[43] ✓ 40.0s  
...  
RandomForestClassifier  
RandomForestClassifier()  
  
Model Validation  
▷  
from sklearn.metrics import accuracy_score, recall_score  
  
pred_y = model.predict(features_validate_x)  
print(accuracy_score(validate_y, pred_y))  
print(recall_score(validate_y, pred_y))  
[44] ✓ 3.0s  
...  
0.8489373516208296  
0.7846326154428514
```

```
from sklearn.ensemble import RandomForestClassifier  
from sklearn.svm import LinearSVC  
  
model = RandomForestClassifier(n_estimators=200)  
model.fit(features_train_x, train_y)  
[41] ✓ 1m 20.3s  
...  
RandomForestClassifier  
RandomForestClassifier(n_estimators=200)  
  
Model Validation  
▷  
from sklearn.metrics import accuracy_score, recall_score  
  
pred_y = model.predict(features_validate_x)  
print(accuracy_score(validate_y, pred_y))  
print(recall_score(validate_y, pred_y))  
[42] ✓ 5.7s  
...  
0.8498471596883845  
0.7847726159887197
```

`RandomForestClassifier`; esto nos dio un accuracy de **0.8489** y un recall de **0.7846** que son buenos indicadores para el modelo ya que están cercanos a 1.

Luego entrenamos nuestro modelo usando 200 estimadores y tomando los demás parámetros por default del algoritmo; esto nos dio un accuracy de **0.8490** y un recall de **0.7847** que no son un incremento significativo de nuestras métricas iniciales y con un tiempo de entrenamiento mayor al inicial.

Ya analizado que 100 estimadores era el parámetro ideal, procedimos a entrenar nuestro modelo usando 100 estimadores, con una profundidad de 10 (# nodos/hojas de profundidad del árbol) y tomando los demás parámetros por default del algoritmo; esto nos dio un accuracy de **0.7485** y un recall de **0.4620** que nos indica un empeoramiento en las métricas.

Con el experimento anterior, decidimos aumentar la profundidad a 100 y esto dio un accuracy de **0.8486** y un recall de **0.7830** que mejoran las métricas y se acercan a las iniciales.

```
from sklearn.ensemble import RandomForestClassifier
from sklearn.svm import LinearSVC

model = RandomForestClassifier(n_estimators=100, max_depth=10)
model.fit(features_train_x, train_y)

[45] ✓ 212s
...
+ RandomForestClassifier
  RandomForestClassifier(max_depth=10)

Model Validation

from sklearn.metrics import accuracy_score, recall_score

pred_y = model.predict(features_validate_x)

print(accuracy_score(validate_y, pred_y))
print(recall_score(validate_y, pred_y))

[46] ✓ 0.6s
...
0.7485757448899729
0.4620482068273561
```

```
from sklearn.ensemble import RandomForestClassifier
from sklearn.svm import LinearSVC

model = RandomForestClassifier(n_estimators=100, max_depth=100)
model.fit(features_train_x, train_y)

[47] ✓ 41.4s
...
+ RandomForestClassifier
  RandomForestClassifier(max_depth=100)

Model Validation

from sklearn.metrics import accuracy_score, recall_score

pred_y = model.predict(features_validate_x)

print(accuracy_score(validate_y, pred_y))
print(recall_score(validate_y, pred_y))

[48] ✓ 3.0s
...
0.8486467301286835
0.783092610308701
```

Parámetros finales

Una vez realizado todos estos experimentos donde realizamos un juego de ejecuciones cambiando los parámetros de estimadores y profundidad, llegamos a la conclusión que entrenar el modelo con 100 estimadores y los demás parámetros en default era suficiente para nuestro modelo de clasificación.

```
[38] final_inference_pipeline.fit(final_training_dataset, final_training_response)
```

```
... Pipeline
      feature_engineering: Pipeline
      features: FeatureUnion
      categorical:
        * one_hot_encode
          * OneHotEncoder
      categorical_binarized:
        * binarizer: ColumnTransformer
          * binarizer
            * Binarizer
      scaled:
        * scaler
          * RobustScaler
      pass:
        * passthrough
          * passthrough
      RandomForestClassifier
```

Model testing

```
> test_pred_y = final_inference_pipeline.predict(test_x)
print(accuracy_score(test_pred_y, test_y))
print(recall_score(test_pred_y, test_y))
```

```
[39] ... 0.8512754707365241
      0.8234125542069434
```