

# Análisis y Reporte sobre el desempeño del modelo

Leonardo Alvarado Menéndez  
A01705998

## Introducción:

Se realizó un modelo de regresión logística para analizar los datos recibidos de imágenes de pasas, y con estos poder clasificar si las pasas eran de tipo Kecimen o Besni.

Los datos usados fueron:

- Area
- Perimeter
- MajorAxisLength
- MinorAxisLength
- Eccentricity
- ConvexArea
- Extent

Para mayor información o sobre los datos, puede consultar en la siguiente página:

<https://archive.ics.uci.edu/ml/datasets/Raisin+Dataset#>

## Método empleado

El modelo de regresión logística empleado, para este análisis fue LogisticRegression de Sklearn[1], donde usando como hyperparameters:

- `penalty = 'l1'`
- `solver = 'saga'`

Y corriendo nuestro modelo, separando el dataset como train y test, obtenemos un accuracy de **0.5111** y un MSE de **0.4888** en nuestros datos test y un accuracy de **0.4962** y un MSE de **0.5037** en los datos de train.

En la Fig 1, se pueden ver los resultados obtenidos en el test:

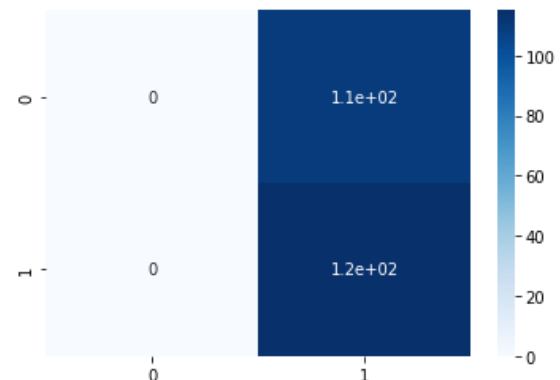


Fig 1: Matriz de confusión modelo1

En esta matriz, el 1 representa las pasas Bensí y el 2 a las Kecimen, como podemos ver, el modelo predice que todas las pasas son de tipo 1 (Bensí), cuando realmente solo 120 de estas son de este tipo y las 110 restantes son de tipo 2 (Kecimen).

Con esta información, podemos ver que nuestro modelo está “**underfitting**”, ya que tenemos un error muy alto tanto en train como en test, por lo que tenemos una bias alta y una varianza baja, ya que nuestro modelo siempre apunta al mismo resultado ‘1’ y únicamente termina el train ya que tiene una limitante de 100 epochs.

## Mejora

En los casos de underfitting, para mejorar nuestro modelo tenemos las siguientes opciones:

1. Transformamos los datos
2. Agregamos más parámetros
3. Modelo más complejo

Para mejorar este modelo se modificó el hyperparameter de solver, el cual originalmente era ‘**saga**’, este hyperparameter, es el solucionador de elección para la regresión logística

1 [https://scikit-learn.org/stable/modules/generated/sklearn.linear\\_model.LogisticRegression.html](https://scikit-learn.org/stable/modules/generated/sklearn.linear_model.LogisticRegression.html)

multinomial dispersa y también es adecuado para conjuntos de datos muy grandes, pero en nuestro caso al no tener un dataset de gran tamaño nos genera un modelo de baja fidelidad, para mejorar esto, podemos usar el solver **'liblinear'**, el cual, aunque es más lento en datasets grandes, es mejor para datasets más pequeños y es la segunda opción cuando usamos como penalty en el modelo **'l1'**.

una precisión de más de 90% ya genera una mejor predicción.

Con estos cambios, volvemos a correr nuestro modelo, y en esta ocasión obtenemos un accuracy de **0.8933** y un MSE de **0.1066** en nuestros datos test y un accuracy de **0.8577** y un MSE de **0.1422** en los datos de train.

En la Fig 2, se pueden ver los resultados obtenidos:

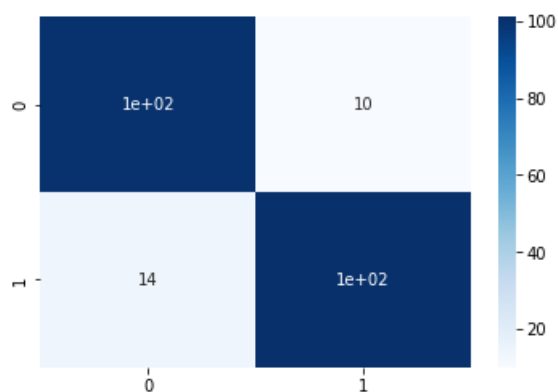


Fig 2:Matriz de confusión modelo 2

En esta ocasión, podemos ver que el modelo únicamente predice mal 10 pasas que llama como tipo 1 cuando realmente son tipo 0 y 14 como tipo 0 cuando son tipo 1.

## Conclusión

Comparando ambos modelos, podemos ver una mejora del **174.78%** en accuracy y una disminución en el error de **0.5111** a **0.1066**, por lo que contamos con un mejor modelo, el cual, aunque no llega a tener