

Universidad de Buenos Aires

FACULTAD DE INGENIERÍA

2DO CUATRIMESTRE 2021

[75.06/95.58]

ORGANIZACIÓN DE DATOS

Trabajo Práctico - Parte 2 Algoritmos de Machine Learning

Integrantes:

Flouret, Agustín Miguel Milhas, Facundo

Mail:

aflouret@fi.uba.ar fmilhas@fi.uba.ar Padrón:

 $102298 \\ 102727$

${\rm \acute{I}ndice}$

1.	Introducción	2
2.	Preprocesamientos utilizados	2
3.	Modelos	2
4.	Conclusión	2

1. Introducción

En este trabajo práctico se utilizaron distintos algoritmos de machine learning para predecir la lluvia de hamburguesas al día siguiente, utilizando diversas técnicas de preprocesamiento y comparando sus métricas.

2. Preprocesamientos utilizados

Nombre del preproce-	Explicación simple	Nombre de la función de Python
samiento		
Básico	Es el preprocesamiento básico utilizado en todos los	basic_preprocessing
	modelos. Esta función se encarga de que haya cohe-	
	rencia en las unidades de los datos, elimina filas en	
	las que la variable target se encuentra vacía, codifica	
	la variable taget de forma binaria, divide el dataset	
	en train y test, y completa valores nulos.	
KNNStandard	Elimina features categoricos y aplica StandardScaler	preprocessing_knn_standard
KNNMinMax	Elimina features categoricos y aplica MinMaxScaler	preprocessing_knn_min_max
KNNNormalizer	Elimina features categoricos y aplica Normalizer	preprocessing_knn_normalizer
Arboles1	Aplica Dummy Encoding a todas las features cate-	preprocessing_arboles_1
	goricas	
Arboles2	Elimina las features categoricas	preprocessing_arboles_2
Arboles3	Aplica Dummy Encoding y elimina columnas con al-	preprocessing_arboles_3
	ta cantidad de nulos	
Arboles4	Elimina features irrelevantes	preprocessing_arboles_4
Redes1	Aplica One Hot Encoding en variables categoricas y	preprocessing_redes_1
	Standard Scaler en variables numericas	
Redes2	Elimina features irrelevantes y aplica Standard Sca-	preprocessing_redes_2
	ler en variables numericas	

3. Modelos

Nombre Modelo	Nombre del preprocesamiento	AUC-ROC	Accuracy	Precision	Recall	F1 score
Arboles de Decisión	Arboles1	0.86	0.85	0.72	0.49	0.58
Naive Bayes	Arboles2	0.83	0.83	0.66	0.46	0.54
Random Forest	Arboles1	0.87	0.85	0.74	0.51	0.61
Boosting	Arboles1	0.88	0.85	0.74	0.52	0.61
KNN	KNNStandard	0.88	0.85	0.78	0.45	0.57
Redes Neuronales	Redes1	0.90	0.86	0.74	0.58	0.65

4. Conclusión

Analizando los resultados de los modelos, podemos concluir que el mejor modelo obtenido en este trabajo es la red neuronal, ya que obtuvo el mayor score AUC-ROC, con un valor de 0.90. Por otro lado, el peor modelo es Naive Bayes, ya que obtuvo el menor puntaje en todas las métricas.

Sin embargo, si se quisiera tener la menor cantidad de falsos positivos, KNN es el más indicado, ya que posee el mayor score de Precision. Por otro lado, si se necesitase obtener la mayor cantidad de dias que llueven hamburguesas sin preocuparse por los falsos positivos, la métrica que hay que maximizar es el Recall. En ese caso, la red neuronal es superior a todos los otros modelos.