

Tercera parte: investigación de modelos supervisados

Random forest classifier

Es un modelo que se utiliza en machine learning para aprendizaje supervisado.

Se usa tanto en problemas de clasificación como en regresión (para predecir valores numéricos).

Random Forest crea muchos árboles de decisión, lo que soluciona el sobreajuste de un solo árbol de decisión que tiende a memorizar los datos de entrenamiento y fallar al ver datos nuevos.

Se asegura de que cada árbol sea diferente a los demás. Si todos los árboles fueran iguales, no ganaríamos nada.

```

1 # TERCERA PARTE: INVESTIGACIÓN DE MODELOS SUPERVISADOS
2
3 # dividimos el conjunto de datos
4 X_train_rf, X_test_rf, y_train_rf, y_test_rf = train_test_split(
5     X, Y,
6     test_size=0.2,
7     stratify=Y,
8     random_state=SEED,
9     shuffle = True
10 )
11
12 # creamos el modelo
13 model_rf = random_forest_model(350, SEED, 1)
14
15 # validación cruzada
16 skf = StratifiedKFold(
17     n_splits=5, # 80/20
18     shuffle=True,
19     random_state=SEED
20 )
21
22 cv_scores = cross_val_score(model_rf, X_train_rf, y_train_rf, cv=skf, scoring='recall_macro')
23
24 print('\n---Resultados de la validación cruzada---')
25 print(f'Puntuación media: {cv_scores.mean()}')
26 print(f'Desviación estándar: {cv_scores.std()}')
27
28 # entrenamiento del modelo
29 model_rf.fit(X_train_rf, y_train_rf)
30
31 # predicción
32 y_pred_rf = model_rf.predict(X_test_rf)
33
34 # evaluación
35 print('\n---Resultados del modelo (Random Forest)---')
36 print(f'Accuracy: {accuracy_score(y_test_rf, y_pred_rf)}')
37 print(classification_report(y_test_rf, y_pred_rf))
38
39 # visualización
40 cm = confusion_matrix(y_test_rf, y_pred_rf, labels=model_rf.classes_)
41 disp = ConfusionMatrixDisplay(confusion_matrix=cm, display_labels=model_rf.classes_)
42 disp.plot(cmap=plt.cm.Purples)
43 plt.title('Matriz de confusión - Random Forest')
44 plt.show()
45
46 importances = model_rf.feature_importances_
47 feature_names = df.columns.drop(4)
48 feature_df = pd.DataFrame({'Feature': feature_names, 'Importance': importances})
49 feature_df = feature_df.sort_values(by='Importance', ascending=True)
50 feature_df.plot(kind='barh', x='Feature', y='Importance', color='purple', legend=False)
51 plt.title('Importancia de las características - Random Forest')
52 plt.show()

```

A partir de este modelo, hemos conseguido los siguientes resultados:

---Resultados de la validación cruzada---

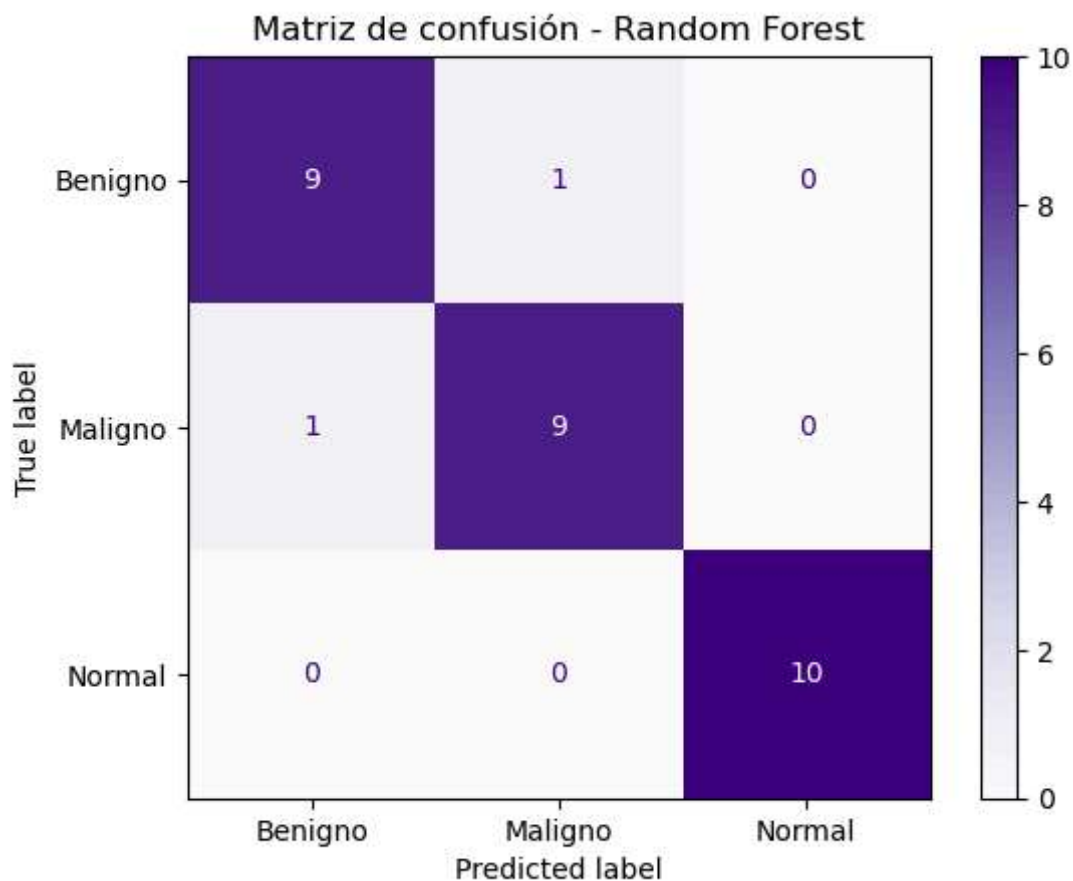
Puntuación media: 0.95

Desviación estándar: 0.0311804782231162

---Resultados del modelo (Random Forest)---

Accuracy: 0.9333333333333333

	precision	recall	f1-score	support
Benigno	0.90	0.90	0.90	10
Maligno	0.90	0.90	0.90	10
Normal	1.00	1.00	1.00	10
accuracy			0.93	30
macro avg	0.93	0.93	0.93	30
weighted avg	0.93	0.93	0.93	30



La siguiente gráfica muestra que variables de las 4 son las que más peso tienen a la hora de hacer una predicción:

