

```

1 # Conjunto de datos
2 import pandas as pd
3 import numpy as np
4 # https://archive.ics.uci.edu/ml/machine-learning-databases/breast-cancer-wisconsin/wdbc.dat
5 df = pd.read_csv('https://bit.ly/3gob0mX', header=None)

```

```
1 df.head(2)
```

```

➡

```

	0	1	2	3	4	5	6	7	8	9	...	22	23	24
0	842302	M	17.99	10.38	122.8	1001.0	0.11840	0.27760	0.3001	0.14710	...	25.38	17.33	184.6
1	842517	M	20.57	17.77	132.9	1326.0	0.08474	0.07864	0.0869	0.07017	...	24.99	23.41	158.8

2 rows x 32 columns

```

1 # Predictoras y objetivo
2 X = df.loc[:, 2:].values
3 y = df.loc[:, 1].values

```

```

1 # Codificación de etiquetas
2 from sklearn.preprocessing import LabelEncoder
3 le = LabelEncoder()
4 y = le.fit_transform(y)
5 le.classes_, le.transform(['M','B']) # clases y ejemplo

```

```
➡ (array(['B', 'M'], dtype=object), array([1, 0]))
```

```

1 # Entrenamiento y pruebas
2 from sklearn.model_selection import train_test_split
3 X_train, X_test, y_train, y_test = train_test_split(X, y, test_size=0.2,
4                                                    stratify=y, random_state=1)
5 X_train.shape, X_test.shape

```

```
➡ ((455, 30), (114, 30))
```

```

1 # Primer pipeline
2 from sklearn.preprocessing import StandardScaler
3 from sklearn.decomposition import PCA
4 from sklearn.linear_model import LogisticRegression
5 from sklearn.pipeline import make_pipeline
6
7 pipe_lr = make_pipeline(StandardScaler(),
8                          PCA(n_components=2),
9                          LogisticRegression(random_state=1))
10 pipe_lr.fit(X_train, y_train)
11 y_pred = pipe_lr.predict(X_test)
12 print('Exactitud en test = %.3f' % pipe_lr.score(X_test, y_test))

```

```
➡ Exactitud en test = 0.956
```

```

1 # Validación cruzada en sklearn
2 from sklearn.model_selection import cross_val_score
3
4 scores = cross_val_score(estimator=pipe_lr,
5                           X=X_train, y=y_train,
6                           cv=10, n_jobs=1) # njobs=-1 => indica usar todos los núcleos disponibles
7 print('Puntajes de exactitud de validación cruzada : %s'%(scores))
8 print('Exactitud de validación cruzada : %.3f +/- %.3f'%(np.mean(scores), np.std(scores)))

```

⇒ Puntajes de exactitud de validación cruzada : [0.93478261 0.93478261 0.95652174 0.95652174 0.97777778 0.93333333 0.95555556 0.95555556]

Exactitud de validación cruzada : 0.950 +/- 0.014

```

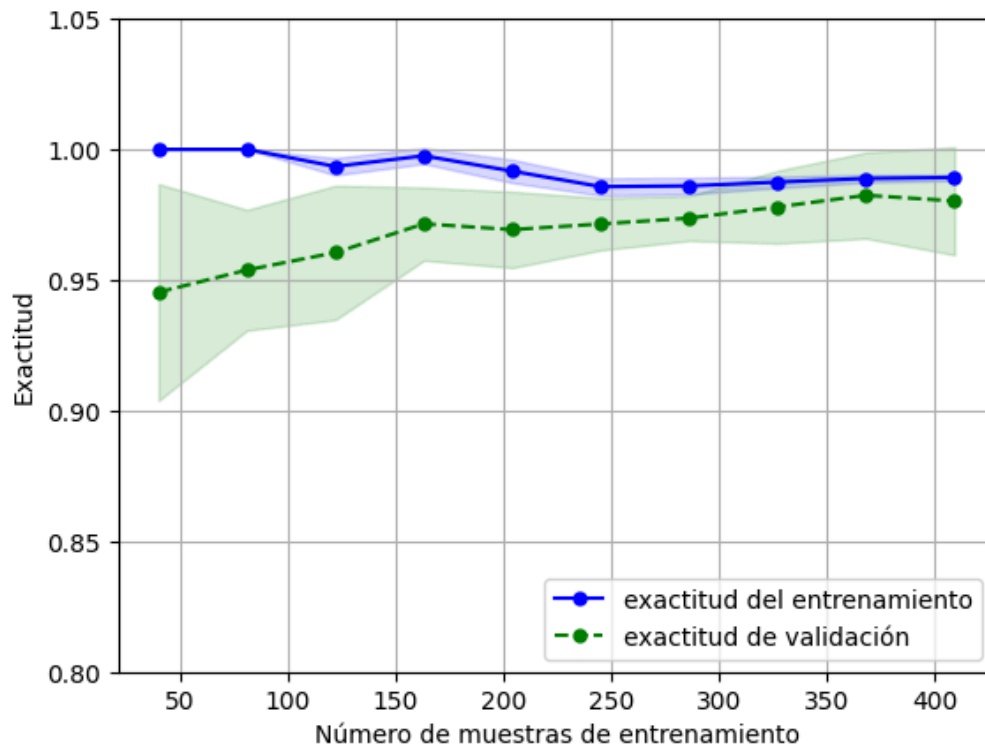
1 # Curvas de aprendizaje
2 import matplotlib.pyplot as plt
3 from sklearn.model_selection import learning_curve
4
5 pipe_lr = make_pipeline(StandardScaler(),
6                           LogisticRegression(penalty='l2', random_state=1))
7
8 train_sizes, train_scores, test_scores = learning_curve(estimator=pipe_lr,
9                                                         X=X_train, y=y_train,
10                                                         train_sizes=np.linspace(0.1, 1, 10),
11                                                         cv=10, n_jobs=1)

```

```

1 train_mean = np.mean(train_scores, axis=1)
2 train_std = np.std(train_scores, axis=1)
3 test_mean = np.mean(test_scores, axis=1)
4 test_std = np.std(test_scores, axis=1)
5
6 plt.plot(train_sizes, train_mean, color='blue', marker='o', markersize=5,
7          label='exactitud del entrenamiento')
8 plt.fill_between(train_sizes, train_mean+train_std, train_mean-train_std,
9                  alpha=0.15, color='blue')
10
11 plt.plot(train_sizes, test_mean, color='green', marker='o', markersize=5,
12          linestyle='--', label='exactitud de validación')
13 plt.fill_between(train_sizes, test_mean+test_std, test_mean-test_std,
14                  alpha=0.15, color='green')
15
16 plt.grid()
17 plt.xlabel('Número de muestras de entrenamiento')
18 plt.ylabel('Exactitud')
19 plt.legend(loc='lower right')
20 plt.ylim([0.8, 1.05])
21 plt.show()

```



```
1 # Sobre/subajuste & curvas de validación
2 from sklearn.model_selection import validation_curve
3 param_range = [10**i for i in range(-3,3)]
4 train_scores, test_scores = validation_curve(estimator=pipe_lr,
5                                             X=X_train, y=y_train,
6                                             param_name='logisticregression__C',
7                                             param_range=param_range,
8                                             cv=10, n_jobs=10)
```

```
1 param_range = [10**i for i in range(-3,3)]
2 param_range
```



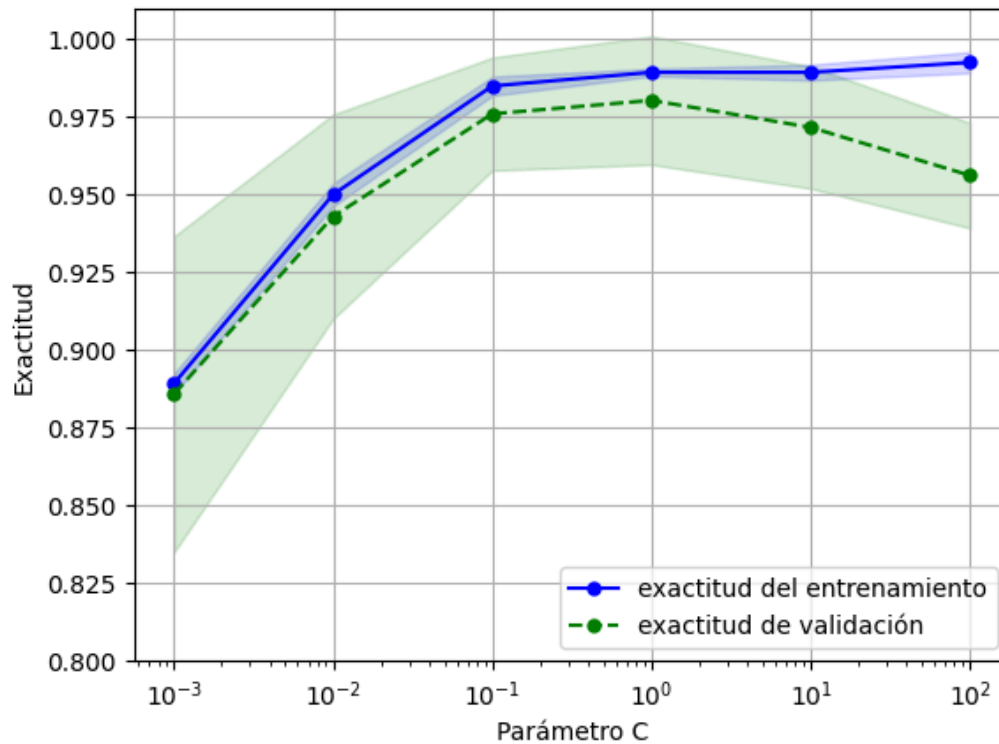
```
[0.001, 0.01, 0.1, 1, 10, 100]
```

```
1 train_mean = np.mean(train_scores, axis=1)
2 train_std = np.std(train_scores, axis=1)
3 test_mean = np.mean(test_scores, axis=1)
4 test_std = np.std(test_scores, axis=1)
5
6 plt.plot(param_range, train_mean, color='blue', marker='o', markersize=5,
7          label='exactitud del entrenamiento')
8 plt.fill_between(param_range, train_mean+train_std, train_mean-train_std,
9                  alpha=0.15, color='blue')
10
11 plt.plot(param_range, test_mean, color='green', marker='o', markersize=5,
12          linestyle='--', label='exactitud de validación')
13 plt.fill_between(param_range, test_mean+test_std, test_mean-test_std,
14                  alpha=0.15, color='green')
15
16 plt.grid()
17 plt.xscale('log')
18 plt.xlabel('Parámetro C')
```

```

19 plt.ylabel('Exactitud')
20 plt.legend(loc='lower right')
21 plt.ylim([0.8, 1.01])
22 plt.show()

```



```

1 # Ajuste de hiperparámetros con búsqueda de malla => grid search
2 from sklearn.model_selection import GridSearchCV
3 from sklearn.svm import SVC
4
5 pipe_svc = make_pipeline(StandardScaler(),
6                           SVC(random_state=1))
7
8 param_range = [10**i for i in range(-4,5)]
9 param_grid_svc = [{'svc__C':param_range,'svc__kernel':['linear']},
10                  {'svc__C':param_range,'svc__gamma':param_range,
11                  'svc__kernel':['rbf']}]
12
13 gs = GridSearchCV(estimator=pipe_svc, param_grid=param_grid_svc,
14                  scoring='accuracy', cv=10, n_jobs=-1)
15
16 gs = gs.fit(X_train, y_train)
17
18 print(gs.best_score_)
19 print(gs.best_params_)

```



```

0.9846859903381642
{'svc__C': 100, 'svc__gamma': 0.001, 'svc__kernel': 'rbf'}

```

```

1 # exactitud del mejor estimador
2 clf = gs.best_estimator_
3 clf.fit(X_train, y_train)
4 print('Exactitud en test : %.3f' % clf.score(X_test, y_test))

```

⇒ Exactitud en test : 0.974

```
1 # Comparación de búsqueda aleatoria y de malla para estimar hiperparámetros
2 # https://scikit-learn.org/stable/auto\_examples/model\_selection/plot\_randomized\_search.html
```

```
1 # Selección de algoritmos con validación cruzada anidada
2 gs_svc = GridSearchCV(estimator=pipe_svc, param_grid=param_grid_svc,
3                       scoring='accuracy', cv=2)
4
5 scores = cross_val_score(gs_svc, X_train, y_train,
6                           scoring='accuracy', cv=5)

1 print('Exactitud NCV con SVC : %.3f +/- %.3f'%(np.mean(scores),np.std(scores)))
```

⇒ Exactitud NCV con SVC : 0.974 +/- 0.015

```
1 from sklearn.tree import DecisionTreeClassifier
2
3 gs_dtc = GridSearchCV(estimator=DecisionTreeClassifier(random_state=0),
4                       param_grid=[{'max_depth': [1,2,3,4,5,6,7,None]}],
5                       scoring='accuracy', cv=2)
6
7 scores = cross_val_score(gs_dtc, X_train, y_train,
8                           scoring='accuracy', cv=5)

1 print('Exactitud de NCV con DT : %.3f +/- %.3f'%(np.mean(scores),np.std(scores)))
```

⇒ Exactitud de NCV con DT : 0.934 +/- 0.016

```
1 # Experimento propio: usar pipe_lr

1 # Métricas de rendimiento
2 # Matriz de confusión
3 from sklearn.metrics import confusion_matrix
4
5 pipe_svc.fit(X_train, y_train)
6 y_pred = pipe_svc.predict(X_test)
7 confmat = confusion_matrix(y_true=y_test, y_pred=y_pred)
8 print(confmat)
```

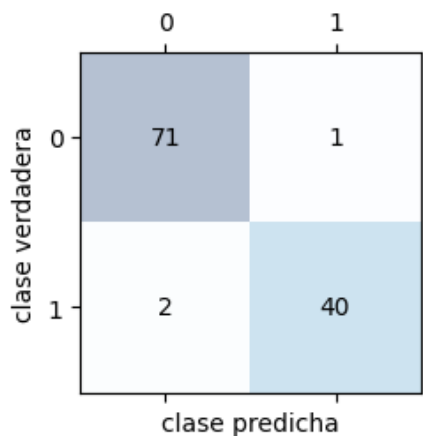
⇒

```
[[71  1]
 [ 2 40]]
```

```

1 confmat = np.array([[71,1],[2,40]])
2 fig, ax = plt.subplots(figsize=(2.5, 2.5))
3 ax.matshow(confmat, cmap=plt.cm.Blues, alpha=0.3)
4 for i in range(confmat.shape[0]):
5     for j in range(confmat.shape[1]):
6         ax.text(x=j, y=i, s=confmat[i,j],
7                 va='center', ha='center')
8 plt.xlabel('clase predicha')
9 plt.ylabel('clase verdadera')
10 plt.show()

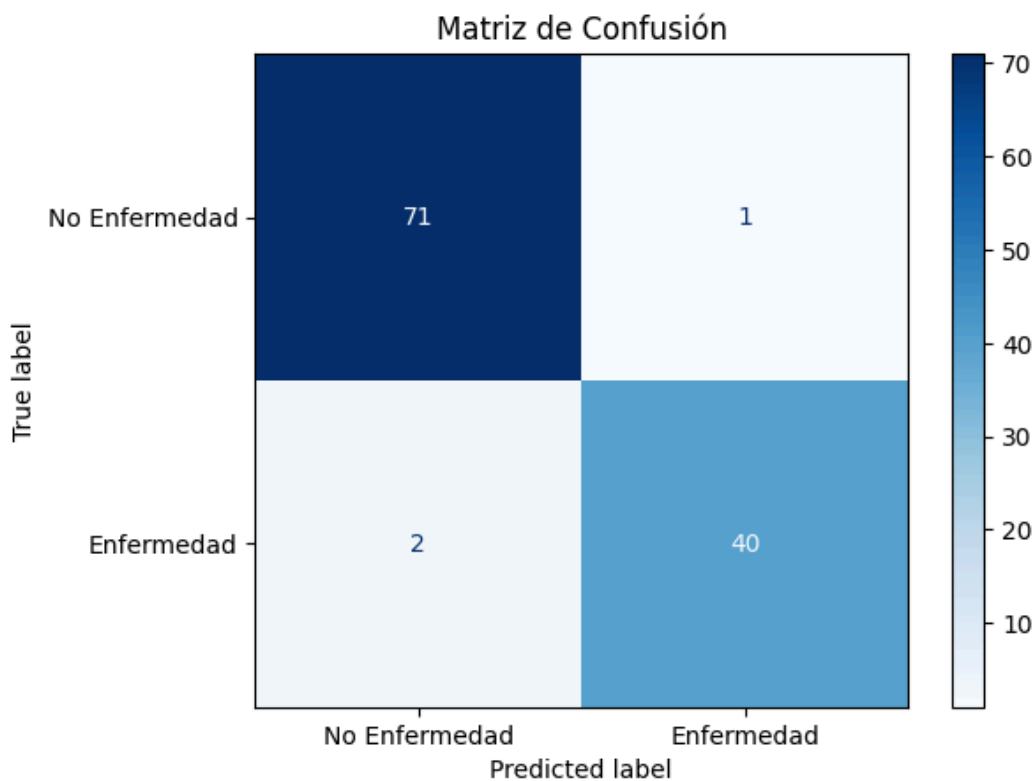
```



```

1 # Con sklearn
2 from sklearn.metrics import ConfusionMatrixDisplay
3 disp = ConfusionMatrixDisplay(confusion_matrix=confmat, display_labels=['No Enfermedad', 'Er
4 disp.plot(cmap=plt.cm.Blues) # Puedes cambiar el mapa de color
5 plt.title('Matriz de Confusión')
6 plt.show()

```



```

1 # Precisión y sensibilidad
2 from sklearn.metrics import precision_score
3 from sklearn.metrics import recall_score, f1_score
4
5 print('Precisión : %.3f' % precision_score(y_true=y_test, y_pred=y_pred))
6 print('    Recall : %.3f' % recall_score(y_true=y_test, y_pred=y_pred))
7 print('        F1 : %.3f' % f1_score(y_true=y_test, y_pred=y_pred))

```

```

⇒ Precisión : 0.976
    Recall : 0.952
    F1 : 0.964

```

```

1 print(gs.best_score_)
2 print(gs.best_params_)

```

```

⇒ 0.9846859903381642
{'svc__C': 100, 'svc__gamma': 0.001, 'svc__kernel': 'rbf'}

```

```

1 # Gráfica de ROC (Receiver Operating Characteristic, o Característica Operativa del Receptor)
2 from sklearn.metrics import roc_curve, auc
3 from numpy import interp
4 from sklearn.model_selection import StratifiedKFold
5
6 pipe_lr = make_pipeline(StandardScaler(),
7                          PCA(n_components=2),
8                          LogisticRegression(penalty='l2', random_state=1,
9                                              C=100.0))
10
11 X_train2 = X_train[:, [4, 14]]
12
13 cv = list(StratifiedKFold(n_splits=3).split(X_train, y_train))

```

```

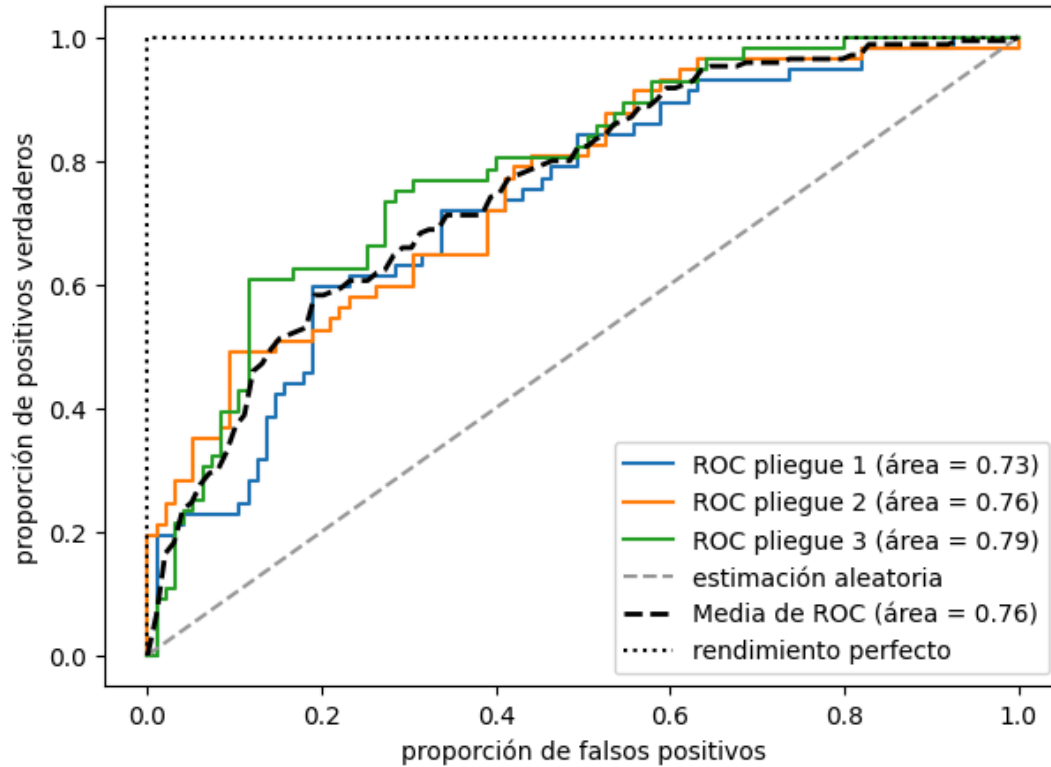
1 fig = plt.figure(figsize=(7, 5))
2
3 mean_tpr = 0.0
4 mean_fpr = np.linspace(0, 1, 100)
5 all_tpr = []
6
7 for i, (train, test) in enumerate(cv):
8     probas = pipe_lr.fit(X_train2[train],
9                         y_train[train]).predict_proba(X_train2[test])
10    fpr, tpr, threshold = roc_curve(y_train[test], probas[:,1], pos_label=1)
11    mean_tpr += interp(mean_fpr, fpr, tpr)
12    mean_tpr[0] = 0.0
13    roc_auc = auc(fpr, tpr)
14    plt.plot(fpr, tpr, label='ROC pliegue %d (área = %0.2f)' % (i+1, roc_auc))
15
16 plt.plot([0,1], [0,1], linestyle='--', color=[0.6,0.6,0.6],
17          label='estimación aleatoria')
18
19 mean_tpr /= len(cv)
20 mean_tpr[-1] = 1.0
21 mean_auc = auc(mean_fpr, mean_tpr)
22 plt.plot(mean_fpr, mean_tpr, 'k--',
23          label='Media de ROC (área = %0.2f)' % mean_auc, lw=2)
24 plt.plot([0,0,1], [0,1,1], linestyle=':', color='black',

```

```

25     label='rendimiento perfecto')
26
27 plt.xlim([-0.05, 1.05])
28 plt.ylim([-0.05, 1.05])
29 plt.xlabel('proporción de falsos positivos')
30 plt.ylabel('proporción de positivos verdaderos')
31 plt.legend(loc='lower right')
32 plt.show()

```



1

1

1

1

1

1

1

1

1

1