

```
In [1]: import numpy as np # arrays and data manipulation
import matplotlib.pyplot as plt # data visualization
import pandas as pd # read, clean, and rearrange data
import seaborn as sns # more styles and colors for data visualization on top of matplotlib

from sklearn.datasets import load_breast_cancer # import dataset from sklearn
from sklearn.model_selection import train_test_split # divide data into train and test parts
from sklearn.preprocessing import StandardScaler # feature scaling so that data works with models better
from sklearn.svm import SVC # support vector machine -> support vector classification: scales quadratically
from sklearn.metrics import classification_report, confusion_matrix
    # classification report gives metrics like accuracy, precision, f1-score, support
    # confusion matrix gives a confusion matrix (TP, TN, FP, FN) between y_true and y_predict

breast_cancer = load_breast_cancer()
```

```
In [2]: breast_cancer
```

[illegible]

```
his is a copy of UCI ML Breast Cancer Wisconsin (Diagnostic) datasets.\nhttps://goo.gl/U2Uwz2\n\nFeatures are computed from a digitized image of a fine needle aspirate (FNA) of a breast mass. They describe characteristics of the cell nuclei present in the image.\n\nSeparating plane described above was obtained using\n\nMultisurface Method-Tree (MSM-T) [K. P. Bennett, "Decision Tree Construction Via Linear Programming." Proceedings of the 4th Midwest Artificial Intelligence and Cognitive Science Society, npp. 97-101, 1992], a classification method which uses linear programming to construct a decision tree. Relevant features were selected using an exhaustive search in the space of 1-4 features and 1-3 separating planes.\n\nThe actual linear program used to obtain the separating plane in the 3-dimensional space is that described in:\n[K. P. Bennett and O. L. Mangasarian: "Robust Linear Programming Discrimination of Two Linearly Inseparable Sets", Optimization Methods and Software 1, 1992, 23-34].\n\nThis database is also available through the UW CS ftp server:\nftp ftp.cs.wisc.edu\ncd math-prog/cpo-dataset/machine-learn/WDBC\n\n.. dropdown:: References\n\n- W.N. Street, W.H. Wolberg and O.L. Mangasarian. Nuclear feature extraction for breast tumor diagnosis. IS&T/SPIE 1993 International Symposium on Electronic Imaging: Science and Technology, volume 1905, pages 861-870, San Jose, CA, 1993.\n- O.L. Mangasarian, W.N. Street and W.H. Wolberg. Breast cancer diagnosis and prognosis via linear programming. Operations Research, 43(4), pages 570-577, July-August 1995.\n- W.H. Wolberg, W.N. Street, and O.L. Mangasarian. Machine learning techniques to diagnose breast cancer from fine-needle aspirates. Cancer Letters 77 (1994) 163-171.\n\n'feature_names': array(['mean radius', 'mean texture', 'mean perimeter', 'mean area', 'mean smoothness', 'mean compactness', 'mean concavity', 'mean concave points', 'mean symmetry', 'mean fractal dimension', 'radius error', 'texture error', 'perimeter error', 'area error', 'smoothness error', 'compactness error', 'concavity error', 'concave points error', 'symmetry error', 'fractal dimension error', 'worst radius', 'worst texture', 'worst perimeter', 'worst area', 'worst smoothness', 'worst compactness', 'worst concavity', 'worst concave points', 'worst symmetry', 'worst fractal dimension'], dtype='<U23'),\n'filename': 'breast_cancer.csv',\n'data_module': 'sklearn.datasets.data'}
```

```
In [4]: df_data = np.c_[breast_cancer['data'], breast_cancer['target']] # stack data array next to target array along 0th axis
df_cols = np.append(breast_cancer['feature_names'], ['target']) # concatenate feature names with a column name
breast_cancer_df = pd.DataFrame(df_data, columns = df_cols)
```

```
In [6]: breast_cancer_df.head()
```

Out[6]:

	mean radius	mean texture	mean perimeter	mean area	mean smoothness	mean compactness	mean concavity	mean concave points	mean symmetry	mean fractal dimension	...	worst texture
0	17.99	10.38	122.80	1001.0	0.11840	0.27760	0.3001	0.14710	0.2419	0.07871	...	17.99
1	20.57	17.77	132.90	1326.0	0.08474	0.07864	0.0869	0.07017	0.1812	0.05667	...	23.44
2	19.69	21.25	130.00	1203.0	0.10960	0.15990	0.1974	0.12790	0.2069	0.05999	...	25.53
3	11.42	20.38	77.58	386.1	0.14250	0.28390	0.2414	0.10520	0.2597	0.09744	...	26.58
4	20.29	14.34	135.10	1297.0	0.10030	0.13280	0.1980	0.10430	0.1809	0.05883	...	16.66

5 rows x 31 columns

```
In [7]: breast_cancer_df.shape
```

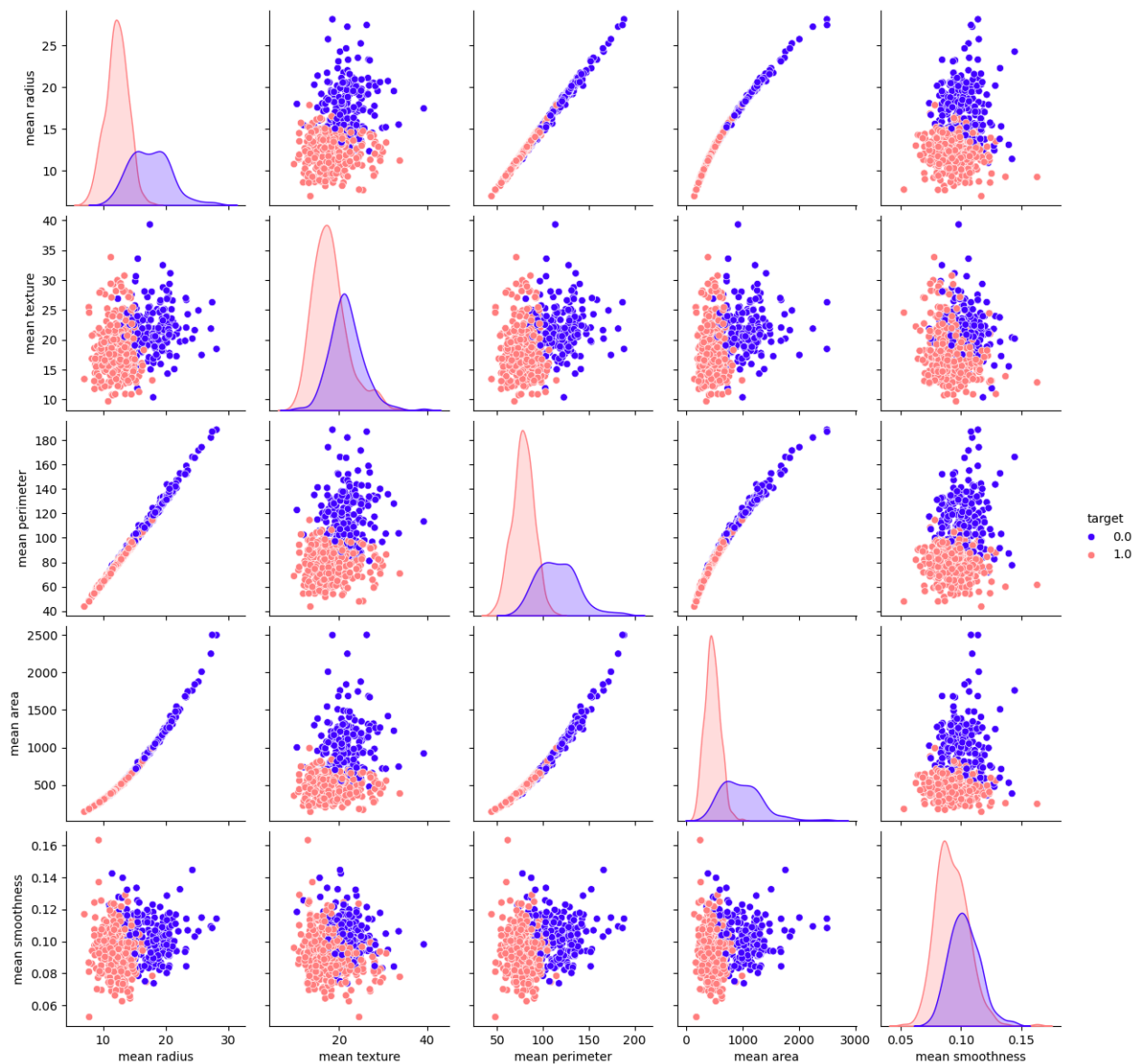
```
Out[7]: (569, 31)
```

```
In [8]: breast_cancer_df.columns
```

```
Out[8]: Index(['mean radius', 'mean texture', 'mean perimeter', 'mean area',
             'mean smoothness', 'mean compactness', 'mean concavity',
             'mean concave points', 'mean symmetry', 'mean fractal dimension',
             'radius error', 'texture error', 'perimeter error', 'area error',
             'smoothness error', 'compactness error', 'concavity error',
             'concave points error', 'symmetry error', 'fractal dimension error',
             'worst radius', 'worst texture', 'worst perimeter', 'worst area',
             'worst smoothness', 'worst compactness', 'worst concavity',
             'worst concave points', 'worst symmetry', 'worst fractal dimension',
             'target'],
            dtype='object')
```

```
In [9]: sns.pairplot(breast_cancer_df, hue = 'target', palette = 'gnuplot2', vars = ['mean radius', 'mean texture',
# 5x5 grid of plots of the different vars plotted against each other
# diagonal plots are KDE of each var (kernel density), non-diagonal plots are scatter plots of one var again
# colors are by target value so red = benign, blue = malignant
```

```
Out[9]: <seaborn.axisgrid.PairGrid at 0x22f54bed9d0>
```

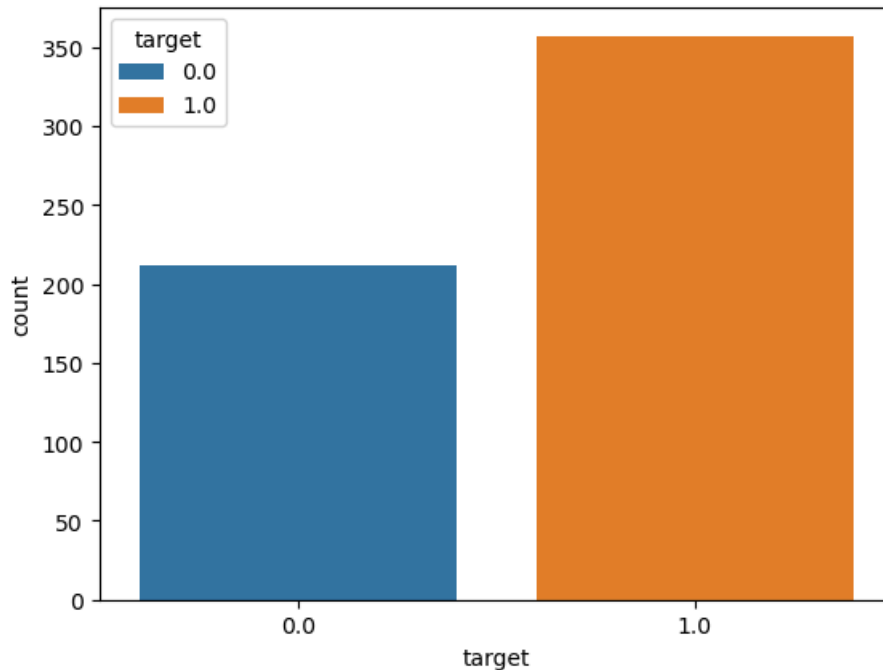


```
In [10]: breast_cancer_df['target'].value_counts() # number of occurrences of each value within targets
```

```
Out[10]: target
1.0      357
0.0      212
Name: count, dtype: int64
```

```
In [11]: sns.countplot(data = breast_cancer_df, x = 'target', hue = 'target') # plot histogram of number of occurrence
```

```
Out[11]: <Axes: xlabel='target', ylabel='count'>
```



```
In [12]: train_data = breast_cancer_df.drop(['target'], axis = 1) # training data removes target values  
train_data.head()
```

```
Out[12]:
```

	mean radius	mean texture	mean perimeter	mean area	mean smoothness	mean compactness	mean concavity	mean concave points	mean symmetry	mean fractal dimension	...	worst radius
0	17.99	10.38	122.80	1001.0	0.11840	0.27760	0.3001	0.14710	0.2419	0.07871	...	25.38
1	20.57	17.77	132.90	1326.0	0.08474	0.07864	0.0869	0.07017	0.1812	0.05667	...	24.99
2	19.69	21.25	130.00	1203.0	0.10960	0.15990	0.1974	0.12790	0.2069	0.05999	...	23.57
3	11.42	20.38	77.58	386.1	0.14250	0.28390	0.2414	0.10520	0.2597	0.09744	...	14.91
4	20.29	14.34	135.10	1297.0	0.10030	0.13280	0.1980	0.10430	0.1809	0.05883	...	22.54

5 rows × 30 columns



```
In [13]: target_data = breast_cancer_df['target'] # target data is only target values  
target_data.head()
```

```
Out[13]:
```

0	0.0
1	0.0
2	0.0
3	0.0
4	0.0

Name: target, dtype: float64

```
In [14]: X_train, X_test, y_train, y_test = train_test_split(train_data, target_data, test_size = 0.2, random_state =  
# split train and target data into training and test datasets, 80% data for training, 20% data for testing,
```

```
In [15]: print("X_train:", X_train.shape)  
print("X_test:", X_test.shape)  
print("y_train:", y_train.shape)  
print("y_test:", y_test.shape)
```

```
X_train: (455, 30)
X_test: (114, 30)
y_train: (455,)
y_test: (114,)
```

```
In [16]: sc = StandardScaler()
X_train = sc.fit_transform(X_train) # fit: compute mean and stdev of X_train, transform: scale with fit data
X_test = sc.transform(X_test) # transform: scale with fit data of X_train
```

```
In [17]: svc_model = SVC()
svc_model.fit(X_train, y_train) # fit the model on training data
```

```
Out[17]: SVC
SVC()
```

```
In [19]: prediction = svc_model.predict(X_test) # predict test data based on fitted model
confusion_mat = np.array(confusion_matrix(y_test, prediction, labels=[1,0])) # create a confusion matrix with
confusion_mat = pd.DataFrame(confusion_mat, index = ["is_malignant", "is_benign"], columns = ["pred_malignant", "pred_benign"])
confusion_mat
```

```
Out[19]:
```

	pred_malignant	pred_benign
is_malignant	66	0
is_benign	1	47

```
In [20]: print(classification_report(y_test, prediction)) # print classification report of pred vs actual data
```

	precision	recall	f1-score	support
0.0	1.00	0.98	0.99	48
1.0	0.99	1.00	0.99	66
accuracy			0.99	114
macro avg	0.99	0.99	0.99	114
weighted avg	0.99	0.99	0.99	114

```
In [ ]:
```