

Introduction to Scikit-Learn (sklearn)

This notebook demonstrates some of the most useful functions of the beautiful Scikit-Learn library.

What we're going to cover:

1. An end-to-end Scikit-Learn workflow
2. Getting the data ready
3. Choose the right estimator/algorithm for our problems
4. Fit the model/algorithm and use it to make predictions on our data
5. Evaluating a model
6. Improve a model
7. Save and load a trained model
8. Putting it all together!

```
In [1]: import pandas as pd
import numpy as np
import matplotlib.pyplot as plt
```

0. An end-to-end Scikit-learn workflow

```
In [2]: # 1. Get the data ready
heart_disease = pd.read_csv("heart-disease.csv")
heart_disease
```

	age	sex	cp	trestbps	chol	fbs	restecg	thalach	exang	oldpeak	slope	ca	thal	target
0	63	1	3	145	233	1	0	150	0	2.3	0	0	1	1
1	37	1	2	130	250	0	1	187	0	3.5	0	0	2	1
2	41	0	1	130	204	0	0	172	0	1.4	2	0	2	1
3	56	1	1	120	236	0	1	178	0	0.8	2	0	2	1
4	57	0	0	120	354	0	1	163	1	0.6	2	0	2	1
...
298	57	0	0	140	241	0	1	123	1	0.2	1	0	3	0
299	45	1	3	140	264	0	1	132	0	1.2	1	0	3	0
300	68	1	0	144	193	1	1	141	0	3.4	1	2	3	0
301	57	1	0	130	151	0	1	115	1	1.2	1	1	3	0
302	57	0	1	130	236	0	0	174	0	0.0	1	1	2	0

303 rows × 14 columns

```
In [3]: # Create the features matrix
X = heart_disease.drop("target", axis=1)
# Create labels
y = heart_disease["target"]
```

```
In [4]: # 3. Choose the right model and hyperparameters
from sklearn.ensemble import RandomForestClassifier

clf = RandomForestClassifier()

# We'll keep the default hyperparameters for now
clf.get_params()
```

```
Out[4]: {'bootstrap': True,
'criterion': 'gini',
'max_depth': None,
'max_features': 'auto',
'max_leaf_nodes': None,
'max_samples': 1,
'min_impurity_decrease': 0.0,
'min_impurity_split': None,
'min_samples_leaf': 1,
'min_samples_split': 2,
'min_weight_fraction_leaf': 0.0,
'n_estimators': 100,
'n_jobs': None,
'oob_score': False,
'random_state': None,
'verbose': 0,
'warm_start': False}
```

```
In [5]: # 3. Fit the model to the training data
from sklearn.model_selection import train_test_split

X_train, X_test, y_train, y_test = train_test_split(X, y,
                                                    test_size=0.2) # 80% training, 20% test
```

```
In [6]: clf.fit(X_train, y_train);
```

```
In [7]: # make a prediction
y_preds = clf.predict(X_test)
y_preds
```

```
Out[7]: array([0, 1, 0, 1, 0, 0, 1, 1, 1, 1, 0, 1, 1, 0, 1, 0, 1, 0, 0, 1, 0, 1, 0, 1, 1, 0,
 0, 1, 0, 1, 0, 1, 1, 1, 1, 1, 1, 1, 0, 0, 1, 1, 1, 0, 1, 1, 0, 0, 1, 0, 0, 1, 0, 1, 0,
 1, 1, 1, 1, 0, 1, 1, 0, 1, 1, 0, 0, 1, 0, 0, 1, 0, 0, 1, 0], dtype=int64)
```

```
In [8]: # 4. Evaluate the model on the training data...
clf.score(X_train, y_train)
```

```
# its 100% accurate because during training the model
# had the chance to see the actual label to try to
# predict the outcome
```

```
Out[8]: 1.0
```

```
In [9]: # ... and test data
clf.score(X_test, y_test)
```

```
Out[9]: 0.868524590163934
```

```
In [10]: # Use some other evaluation metrics
from sklearn.metrics import classification_report, confusion_matrix, accuracy_score
```

```
In [11]: print(classification_report(y_test, y_preds))
```

	precision	recall	f1-score	support
0	0.92	0.79	0.85	29
1	0.83	0.94	0.88	32
accuracy			0.87	61
macro avg	0.88	0.87	0.87	61
weighted avg	0.87	0.87	0.87	61

```
In [12]: confusion_matrix(y_test, y_preds)
```

```
Out[12]: array([[23, 6],
 [ 2, 30]], dtype=int64)
```

```
In [13]: accuracy_score(y_test, y_preds)
```

```
Out[13]: 0.868524590163934
```

```
In [14]: # 5. Improve the model
# Try different amount of n_estimators
```

```
np.random.seed(42)
for i in range(10, 100, 10):
    print(f"Trying model with {i} estimators")
    clf = RandomForestClassifier(i).fit(X_train, y_train)
    print(f"Model accuracy on test set: (clf.score(X_test, y_test) * 100:.2f)%")
    print("")
```

Trying model with 10 estimators
Model accuracy on test set: 85.25%

Trying model with 20 estimators
Model accuracy on test set: 80.33%

Trying model with 30 estimators
Model accuracy on test set: 85.25%

Trying model with 40 estimators
Model accuracy on test set: 85.25%

Trying model with 50 estimators
Model accuracy on test set: 86.92%

Trying model with 60 estimators
Model accuracy on test set: 86.92%

Trying model with 70 estimators
Model accuracy on test set: 86.89%

Trying model with 80 estimators
Model accuracy on test set: 85.25%

Trying model with 90 estimators
Model accuracy on test set: 85.25%

```
In [15]: # 6. Save a model and load it
import pickle
```

```
pickle.dump(clf, open("random_forest_ml.pkl", "wb"))
```

```
In [16]: loaded_model = pickle.load(open("random_forest_ml.pkl", "rb"))
loaded_model.score(X_test, y_test)
```

```
Out[16]: 0.8524590163934426
```

1. Getting the data ready

Three main things we have to do:

1. Split the data into features and labels (usually X and y)
 - then split them into train set and test set
2. Converting non-numerical values to numerical values (feature encoding)
3. Filling (imputing) or disregarding missing values

```
In [17]: heart_disease.head()
```

	age	sex	cp	trestbps	chol	fbs	restecg	thalach	exang	oldpeak	slope	ca	thal	target
0	63	1	3	145	233	1	0	150	0	2.3	0	0	1	1
1	37	1	2	130	250	0	1	187	0	3.5	0	0	2	1
2	41	0	1	130	204	0	0	172	0	1.4	2	0	2	1
3	56	1	1	120	236	0	1	178	0	0.8	2	0	2	1
4	57	0	0	120	354	0	1	163	1	0.6	2	0	2	1

1.1 Split the data into features and labels

```
In [18]: X = heart_disease.drop("target", axis=1) # features
y = heart_disease["target"] # labels
```

Split features and labels into training and test sets

```
In [19]: # Split the data into training and test sets
from sklearn.model_selection import train_test_split

test_size = 0.2
X_train, X_test, y_train, y_test = train_test_split(X, y, test_size=test_size)
```

```
In [20]: X_train.shape, X_test.shape, y_train.shape, y_test.shape
# X_train: 80% of the dataset
# X_test: number of labels (columns)
```

```
Out[20]: (242, 13), (61, 13), (242,), (61,)
```

1.2 Converting non-numerical values

For this, we need another dataset

```
In [21]: car_sales = pd.read_csv("car-sales-extended.csv")
car_sales.head()
```

	Make	Colour	Odometer (KM)	Doors	Price
0	Honda	White	35431	4	15323.0
1	BMW	Blue	192714	5	19943.0
2	Honda	White	84714	4	28343.0
3	Toyota	White	154365	4	13434.0
4	Nissan	Blue	181577	3	14043.0

```
In [22]: car_sales.dtypes
```

```
Out[22]: Make          object
Colour         object
Odometer (KM)  int64
Doors          int64
Price         int64
dtype: object
```

```
In [23]: # Split into X/y
X = car_sales.drop("Price", axis=1)
y = car_sales["Price"]

# Split into training and testing
X_train, X_test, y_train, y_test = train_test_split(X, y, test_size=0.2)
```

```
In [24]: # Build machine learning model
from sklearn.ensemble import RandomForestRegressor

try:
    model = RandomForestRegressor()
    model.fit(X_train, y_train)
    model.score(X_test, y_test)
except ValueError as e:
    print(f"ValueError: {e}")
```

ValueError: could not convert string to float: 'Toyota'

ValueError: could not convert string to float

```
In [25]: # Turn the categories into numbers
from sklearn.preprocessing import OneHotEncoder
from sklearn.compose import ColumnTransformer
```

```
categorical_features = ["Make", "Colour",
                        "Doors"] # doors too, because each car falls into
                                # a "category" of doors (3, 4, 5)

one_hot = OneHotEncoder()
transformer = ColumnTransformer([("one_hot",
                                categorical_features)],
                                remainder="passthrough")
transformed_X = transformer.fit_transform(X)
```

```
Out[25]: array([[0.00000e+00, 1.00000e+00, 0.00000e+00, ..., 1.00000e+00,
 1.00000e+00, 0.00000e+00, 3.54310e+04],
 [0.00000e+00, 0.00000e+00, 0.00000e+00, ..., 0.00000e+00,
 1.00000e+00, 1.92714e+05],
 [0.00000e+00, 1.00000e+00, 0.00000e+00, ..., 1.00000e+00,
 0.00000e+00, 8.47140e+04],
 ...,
 [0.00000e+00, 0.00000e+00, 1.00000e+00, ..., 1.00000e+00,
 0.00000e+00, 6.66040e+04],
 [0.00000e+00, 1.00000e+00, 0.00000e+00, ..., 1.00000e+00,
 0.00000e+00, 2.15883e+05],
 [0.00000e+00, 0.00000e+00, 0.00000e+00, ..., 1.00000e+00,
 0.00000e+00, 2.48360e+05]])
```

```
In [26]: X.head()
```

	Make	Colour	Odometer (KM)	Doors
0	Honda	White	35431	4
1	BMW	Blue	192714	5
2	Honda	White	84714	4
3	Toyota	White	154365	4
4	Nissan	Blue	181577	3

```
In [27]: pd.DataFrame(transformed_X).head()
```

	0	1	2	3	4	5	6	7	8	9	10	11	12
0	0.0	1.0	0.0	0.0	0.0	0.0	0.0	0.0	1.0	0.0	1.0	0.0	35431.0
1	1.0	0.0	0.0	0.0	1.0	0.0	0.0	0.0	0.0	0.0	1.0	0.0	192714.0
2	0.0	1.0	0.0	0.0	0.0	0.0	0.0	0.0	1.0	0.0	1.0	0.0	84714.0
3	0.0	0.0	0.0	1.0	0.0	0.0	0.0	1.0	0.0	1.0	0.0	0.0	154365.0
4	0.0	0.0	1.0	0.0	0.0	1.0	0.0	0.0	1.0	0.0	0.0	0.0	181577.0

What happened here?

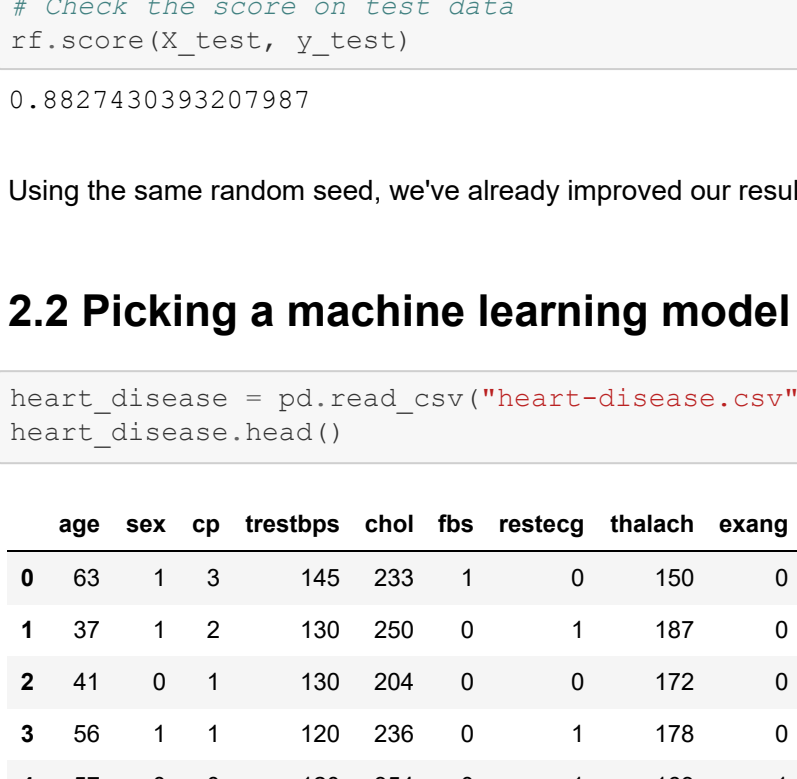
Using the `OneHotEncoder`, we've transformed all the categorical features into numerical features so that our Machine Learning Model can work with them.

How did it do that?

Say that the feature `Make` has three possibilities: `Honda`, `Toyota` and `BMW`. The `OneHotEncoder`

1. transforms the feature `Make` into a vector of size 3 filled with zeros.
 - Say that the first zero represents `Honda`, the second `Toyota` and the third `BMW`
2. for each entry (row), it changes 0 into 1 in the slot that corresponds to the actual categorical feature.
 - Say that the entry was a `Toyota`, then the vector would be `[0,1,0]`

```
In [28]: from IPython.display import YouTubeVideo
YouTubeVideo('v_4f0mKmsU', start=59)
```



```
In [29]: # a better view using a similar method from pandas
dummies = pd.get_dummies(X)
dummies.head()
```

	Odometer (KM)	Doors	Make_BMW	Make_Honda	Make_Nissan	Make_Toyota	Colour_Blue	Colour_Green	Colour_Red	Colour_White
0	35431	4	0	1	0	0	0	0	0	0
1	192714	5	1	0	0	0	0	0	1	0
2	84714	4	0	1	0	0	0	0	0	0
3	154365	4	0	0	0	1	0	0	0	0
4	181577	3	0	0	1	0	0	1	0	0

```
In [30]: # Let's refit the model
np.random.seed(42) # needed if we want the same outcome every time we run this cell
X_train, X_test, y_train, y_test = train_test_split(transformed_X, y, test_size=0.2)
model.fit(X_train, y_train)
model.score(X_test, y_test)
```

```
Out[30]: 0.3235867221569877
```

1.3 Missing values

We have two options for handling missing data:

1. Fill the missing data (imputation)
2. Remove the samples with missing data

Let's import another dataset that has missing values.

```
In [31]: car_sales_missing = pd.read_csv("car-sales-extended-missing-data.csv")
car_sales_missing.head()
```

	Make	Colour	Odometer (KM)	Doors	Price
0	Honda	White	35431.0	4.0	15323.0
1	BMW	Blue	192714.0	5.0	19943.0
2	Honda	White	84714.0	4.0	28343.0
3	Toyota	White	154365.0	4.0	13434.0
4	Nissan	Blue	181577.0	3.0	14043.0

From the head of the `DataFrame`, it doesn't seem like we have missing values.

We can calculate how many there are in this way.

```
In [32]: car_sales_missing.isna().sum()
```

```
Out[32]: Make          49
Colour         50
Odometer (KM)  50
Doors          50
Price         50
dtype: int64
```

Option 1: Fill missing data with Pandas

```
In [33]: # Fill the "Make" column
car_sales_missing["Make"].fillna("missing", inplace=True)
```

```
# Fill the "Colour" column
car_sales_missing["Colour"].fillna("missing", inplace=True)
```

```
# Fill the "Odometer" column
car_sales_missing["Odometer (KM)"].fillna(car_sales_missing["Odometer (KM)"].mean(),
                                           inplace=True)
```

```
# Fill the "Doors" columns
car_sales_missing["Doors"].fillna(car_sales_missing["Doors"].mode().iat[0],
                                  inplace=True)
```

```
In [34]: # Verify that our code worked
car_sales_missing.isna().sum()
```

```
Out[34]: Make          0
Colour         0
Odometer (KM)  0
Doors          0
Price         50
dtype: int64
```

Option 2: Fill missing data with Scikit-learn

```
In [35]: car_sales_missing2 = pd.read_csv("car-sales-extended-missing-data.csv")
car_sales_missing2.head()
```

	Make	Colour	Odometer (KM)	Doors	Price
0	Honda	White	35431.0	4.0	15323.0
1	BMW	Blue	192714.0	5.0	19943.0
2	Honda	White	84714.0	4.0	28343.0
3	Toyota	White	154365.0	4.0	13434.0
4	Nissan	Blue	181577.0	3.0	14043.0

```
In [36]: car_sales_missing2.isna().sum()
```

```
Out[36]: Make          0
Colour         0
Odometer (KM)  0
Doors          0
Price         50
dtype: int64
```

```
In [37]: # We are not able to train our model if some rows
# don't have the "Price" value (our label), so we
# just drop them
car_sales_missing2.dropna(subset=["Price"], inplace=True)
```

```
In [38]: # Split into X & y
X = car_sales_missing2.drop("Price", axis=1)
y = car_sales_missing2["Price"]
```

```
In [39]: # Fill missing values with Scikit-Learn
from sklearn.impute import SimpleImputer
from sklearn.compose import ColumnTransformer
```

```
# Fill categorical values with "missing" and
# numerical values with mean
cat_imp = SimpleImputer(strategy="constant", fill_value="missing")
num_imp = SimpleImputer(strategy="constant", fill_value=4)
door_imp = SimpleImputer(strategy="mean")
```

```
# Define columns
cat_features = ["Make", "Colour"]
door_features = ["Doors"]
num_features = ["Odometer (KM)"]
```

```
# Create an imputer (something that fills missing data)
imputer = ColumnTransformer([("cat_imputer", cat_features),
                              ("cat_imputer", door_features),
                              ("num_imputer", num_features)],
                              remainder="passthrough")
```

```
# Transform the data
filled_X = imputer.fit_transform(X)
```

```
Out[39]: array([[0,'Honda', 'White', 4.0, 35431.0],
 [1,'BMW', 'Blue', 5.0, 192714.0],
 [2,'Honda', 'White', 4.0, 84714.0],
 [3,'Nissan', 'White', 4.0, 66604.0],
 [4,'Honda', 'Blue', 4.0, 215883.0],
 [5,'Toyota', 'Blue', 4.0, 248360.0]], dtype=object)
```

```
In [40]: pd.DataFrame(filled_X,
                    columns=["Make", "Colour", "Odometer (KM)", "Doors"],
                    ).isna().sum()
```

```
Out[40]: Make          0
Colour         0
Odometer (KM)  0
Doors          0
dtype: int64
```

Option 3: Drop missing data

```
In [41]: car_sales_missing.dropna(inplace=True)
```

2. Choosing the right estimator/algorithm for our problem

Scikit-learn uses "estimator" as another term for "machine learning model" or "algorithm".

There are at least two big kinds of algorithms:

- **"Regression"** - predicting a number
- **"Classification"** - predicting whether a sample is one thing or another (healthy, ill)

