



## Introduction to Machine Learning with scikit-learn

### Linear regression

#### Introduction

The objective of this exercise is to become familiar with linear regression. Linear regression was one of the first predictive models to be studied and is today one of the most popular models for practical applications thanks to its simplicity.

#### Univariate Linear Regression

In the univariate linear model, we have two variables:  $y$  called the **target** variable and  $x$  called the **explanatory variable**. Linear regression consists in modeling the link between these two variables by an **affine function**. Thus, the formula of the univariate linear model is given by:

$$y \approx \beta_1 x + \beta_0$$

where:

- $y$  is the variable we want to predict.
- $x$  is the explanatory variable.
- $\beta_1$  and  $\beta_0$  are the parameters of the affine function.  $\beta_1$  will define its **slope** and  $\beta_0$  will define its **y-intercept** (also called **bias**).

The goal of linear regression is to estimate the best parameters  $\beta_0$  and  $\beta_1$  to predict the variable  $y$  from a given value of  $x$ .

To get a feel for Univariate Linear Regression, let us look at the interactive example below.

- (a) Run the next cell to display the interactive figure. In this figure, we have simulated a dataset by the relation  $y = \alpha_1 x + \alpha_0$ .
- (b) Use the sliders on the `Regression` tab to find the parameters  $\beta_0$  and  $\beta_1$  that **best match** all the points in the data set.
- (c) What is the effect of each of the parameters on the regression function?

In [ ]:

```
from widgets import regression_widget  
regression_widget()
```

#### Multivariate Linear Regression

Multivariate linear regression consists in modeling a linear link between a target variable  $y$  and several **explanatory variables**  $x_1, x_2, \dots, x_p$ , often called *features*:

$$\begin{aligned} y &\approx \beta_0 + \beta_1 x_1 + \beta_2 x_2 + \dots + \beta_p x_p \\ &\approx \beta_0 + \sum_{j=1}^p \beta_j x_j \end{aligned}$$

There are now  $p + 1$  parameters  $\beta_j$  to find.

### 1. Using scikit-learn for linear regression

We are now going to learn how to use the **scikit-learn** library in order to solve a Machine Learning problem with a **linear regression**.

During the following exercises, the objective will be to predict the **selling price of a car** based on its **characteristics**.

#### Importing the dataset

The dataset that we will use in the following contains many characteristics about different cars from 1985.

For simplicity, only the numeric variables have been kept and the lines containing missing values have been deleted.

- (a) Import the `pandas` module under the alias `pd`.
- (b) In a `DataFrame` named `df`, import the `automobiles.csv` dataset using the `read_csv` function of `pandas`. This file is located in the same folder as the runtime environment of the notebook.
- (c) Display the first 5 lines of `df` to check if the import was successful.

In [ ]:

```
# Insert your code here
```

Show solution

- The `symboling` variable corresponds to the degree of risk with respect to the insurer (risk of accident, breakdown, etc.).
- The `normalized_losses` variable is the relative average cost per year of vehicle insurance. This value is normalized with respect to cars of the same type (SUV, utility, sports, etc.).
- The following 13 variables concern the technical characteristics of the cars such as width, length, engine displacement, horsepower, etc ...
- The last variable `price` corresponds to the selling price of the vehicle. This is the variable that we will try to predict.

### Separation of the explanatory variables from the target variable

We are now going to create two `DataFrames`, one containing the explanatory variables and another containing the target variable `price`.

In [ ]:

```
# Insert your code here
```

Show solution

### Splitting of the data into training and test sets

We are now going to split our dataset into two sets: A **training** set and a **test** set. This step is **extremely** important when doing Machine Learning.

Indeed, as their names indicate:

- The training set is used to train the model, ie to find the optimal  $\beta_0, \dots, \beta_p$  parameters for this dataset.
- The test set is used to test the trained model by evaluating its ability to **generalize** its predictions on data that it has **never seen**.

A very useful function for doing this is the `train_test_split` function of the `model_selection` submodule of **scikit-learn**.

- (f) Run the following cell to import the `train_test_split` function.

In [ ]:

```
from sklearn.model_selection import train_test_split
```

This function is used as follows:

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

- `X_train` and `y_train` are the explanatory and target variables of the **training** dataset.
- `X_test` and `y_test` are the explanatory and target variables of the **test** dataset.
- The `test_size` parameter corresponds to the **proportion** of the dataset that we want to keep for the test set. In the previous example, this proportion corresponds to 20% of the initial dataset.

- (g) Using the `train_test_split` function, separate the dataset into a training set (`X_train`, `y_train`) and a test set (`X_test`, `y_test`) so that the test set contains **15% of the initial dataset**.

In [ ]:

```
# Insert your code here
```

Show solution

### Training the regression model

To train a linear regression model on this dataset, we will use the **LinearRegression** class contained in the `linear_model` submodule of **scikit-learn**.

- (h) Run the following cell to import the `LinearRegression` class.

In [ ]:

```
from sklearn.linear_model import LinearRegression
```

The `scikit-learn` API makes it easy to train and evaluate models. All `scikit-learn` model classes have the following two methods:

- **fit** : Train the model on the dataset given as input.
- **predict** : Make a prediction from a set of explanatory variables given as input.

Below is an example of training a model with `scikit-learn`:

In [ ]:

```
# Insert your code here
```

Show solution

## Evaluation of the model's performance

In order to evaluate the **quality of the predictions of the model** obtained thanks to the parameters  $\beta_0, \dots, \beta_j$ , there are several metrics already built in the `scikit-learn` library.

One of the most used metrics for regression is the **Mean Squared Error** (MSE) which is defined under the name of `mean_squared_error` in the `metrics` submodule of `scikit-learn`.

This function consists in calculating the average of the squared distances between the **target variables** and the **predictions obtained** thanks to the regression function.

The following interactive figure shows how this error is calculated according to  $\beta_1$ :

- The **blue** dots represent the **dataset** for which we want to evaluate the quality of the predictions. Usually this is the **test** dataset.
- The **red** line is the regression function configured by  $\beta_1$ . In this example,  $\beta_0$  is set to 0 to simplify the illustration.
- The **green** lines are the **distances** between the **target variable** and the **predictions** obtained thanks to the regression function parameterized by  $\beta_1$ .

**The mean squared error is just the average of these squared distances.**

- (m) Run the next cell to display the interactive figure.
- (n) Using the cursor below the figure, try to find a value of  $\beta_1$  that minimizes the Mean Squared Error. Is this value unique?

In [ ]:

```
from widgets import interactive_MSE
interactive_MSE()
```

The `mean_squared_error` function of `scikit-learn` is used as follows:

```
mean_squared_error(y_true, y_pred)
```

where:

- `y_true` contains the true values of the target variable.
- `y_pred` contains the values **predicted** by our model for the same explanatory variables.

- (o) Import the `mean_squared_error` function from the `sklearn.metrics` submodule.
- (p) Evaluate the prediction quality of the model on **training data**. Store the result in a variable named `mse_train`.
- (q) Evaluate model prediction quality on **test data**. Store the result in a variable named `mse_test`.
- (r) Why is the MSE higher on the test dataset?

In [ ]:

```
# Insert your code here
```

Show solution

The mean squared error you will find should be around millions on the test data, which can be difficult to interpret.

This is why we are going to use another metric, the **Mean Absolute Error** which is at the same scale as the target variable.

- (s) Import the `mean_absolute_error` function from the `sklearn.metrics` submodule.
- (t) Evaluate the prediction quality on test and training data using the mean absolute error.
- (u) From the `DataFrame` `df`, calculate the average purchase price on all vehicles. Do the model's predictions seem reliable to you?

In [ ]:

```
# Insert your code here
```

Show solution

## 2. Overfitting the data with another regression model

We have just seen that with the `LinearRegression` class of `skikit-learn`, the model was able to learn on the training data and generalize on the test data with an error rate of 20% on average.

In what follows we will create another regression model that **learns very well on training data but generalizes very poorly on test data**: this is called **overfitting**.

For this we will use a Machine Learning model called **Gradient Boosting Regressor** known for its tendency to overfit.

- (a) Run the following cell to import the `GradientBoostingRegressor` class contained in the `ensemble` submodule of `skikit-learn` and instantiate a `GradientBoostingRegressor` model named `gbr`.

In [ ]:

```
from sklearn.ensemble import GradientBoostingRegressor

# These parameters have been chosen to overfit on purpose
# Do not use them in practice
gbr = GradientBoostingRegressor(n_estimators = 1000,
                                max_depth = 10000,
                                max_features = 15,
                                validation_fraction = 0)
```

- (b) Train the model `gbr` using its `fit` method.
- (c) Make predictions on the test and training datasets. Store these predictions in `y_pred_test_gbr` and `y_pred_train_gbr`.

In [ ]:

```
# Insert your code here
```

Show solution

After instantiating our model, training it on the training data and making the predictions, we must then evaluate its performance.

- (d) Calculate the MSE on the **training** data and the **test** data using the `mean_squared_error` function then display the results.
- (e) Calculate the MAE for the **training** data and the **test** data using the `mean_absolute_error` function then display the results.
- (f) After having calculated the average of the `price` column, calculate the **relative error of the model** on the test set.

In [ ]:

```
# Insert your code here
```

Show solution

Here is an example of results that we could obtain with these two models.

For the **linear regression with `LinearRegression`** we had:

- MAE train lr = 1588.131591267774
- MAE test lr = 2105.5002712214014

For the **regression with `GradientBoostingRegressor`** we have:

- MAE train gbr: 27.533333333339847
- MAE test gbr: 1393.013371545563

The mean absolute error obtained on the training set by the `GradientBoostingRegressor` model is only 27.5 against 1588 for the linear regression. The `GradientBoostingRegressor` model is very powerful and is able to learn the training data almost "**by heart**" which explains this difference in performance.

It is for this reason that the performance of the model should be evaluated on the **test** dataset. Indeed, the average absolute error of the `GradientBoostingRegressor` model is 1393, which is **very far** from the performance obtained on the training data.

This is an example of **blatant overfitting**. Even if the performance of the `GradientBoostingRegressor` is superior to that of the linear regression on the test data, you should always **be wary** of too high a performance.

## 3. Going further: Polynomial Regression

In many cases, the relationship between the variables  $x$  and  $y$  is **not linear**. This does not allow us to use linear regression to predict  $y$ . We could then propose a **quadratic model** such as:

$$y = \beta_0 + \beta_1 x + \beta_2 x^2$$

- (a) Run the next cell to display the interactive figure.
- (b) Find the optimal parameters for  $\beta_0$ ,  $\beta_1$  and  $\beta_2$  that best approximate the data on the scatter plot.
- (c) Set  $\beta_2$  to 0 and vary  $\beta_0$  and  $\beta_1$ . Which model do you recognize?

In [ ]:

```
from widgets import polynomial_regression

polynomial_regression()
```

Polynomial regression is equivalent to performing a classical linear regression from **polynomial functions of the explanatory variable** of arbitrary degree. Polynomial regression is much more flexible than classical linear regression because it can approach any type of continuous function.

When we have several explanatory variables, the polynomial variables can also be calculated by products between the explanatory variables. For example, if we have three variables, then the **second-order** polynomial regression model becomes:

$$y \approx \beta_0 + \beta_1 x_1^2 + \beta_2 x_2^2 + \beta_3 x_3^2 + \beta_4 x_1 x_2 + \beta_5 x_2 x_3 + \beta_6 x_1 x_3$$

If we had more explanatory variables or wanted to increase the degree of polynomial regression, the number of explanatory variables **would explode**, which could induce **overfitting**.

- (d) Run the next cell to display the interactive figure.

The scatter plot was generated with the same trend as the previous figure. The red line corresponds to the optimal polynomial regression function obtained on these data.

- (e) Taking into account the scatter plot in the previous figure, find the **degree** of the polynomial regression that best captures the trend of the data.

In [ ]:

```
from widgets import polynomial_regression2
polynomial_regression2()
```

To train a Polynomial regression model with scikit-learn, we must first calculate the **polynomial variables** from the data. This can be done using the **PolynomialFeatures** class of the `preprocessing` submodule:

```
from sklearn.preprocessing import PolynomialFeatures

poly_feature_extractor = PolynomialFeatures(degree = 2)
```

- The **degree** parameter defines the degree of the polynomial features to be calculated.

The `poly_feature_extractor` object is **not a prediction model**. This type of object is called a **Transformer** and it can be used with the following two methods:

- **fit** : does nothing in this case. This method is generally used to calculate the parameters necessary to apply a transformation to the data.
- **transform** : Applies the transformation to the dataset. In this case, the method returns the polynomial features of the dataset.

These two methods can be called **simultaneously** using the **fit\_transform** method. We can compute the polynomial features on `X_train` and `X_test` as follows:

```
X_train_poly = poly_feature_extractor.fit_transform(X_train)

X_test_poly = poly_feature_extractor.transform(X_test)
```

- (g) Import the `PolynomialFeatures` class from the `preprocessing` submodule of `sklearn`.
- (h) Instantiate an object of class `PolynomialFeatures` with the argument **degree = 3** and name it `poly_feature_extractor`.
- (i) Apply the transformation of `poly_feature_extractor` on `X_train` and `X_test` and store the results in `X_train_poly` and `X_test_poly`.

In [ ]:

```
# Insert your code here
```

Show solution

- (j) Train a linear regression model on the data (`X_train_poly`, `y_train`).
- (k) Evaluate its performance on training data and test data (`X_test_poly`, `y_test`). Are we in an overfitting regime?

In [ ]:

```
# Insert your code here
```

Show solution

## Conclusion and recap

In this course, you have been introduced to solving a regression problem with machine learning.



Validate