



TASK

Supervised Learning - Linear Regression

Visit our website

Introduction

WELCOME TO THE SUPERVISED LEARNING - LINEAR REGRESSION TASK!

In this task we will show a simple machine learning algorithm in action. We will learn about regression analysis, a statistical process used to estimate the relationship between some variables. This classic method is used by machine learning experts, as well as data scientists and statisticians.

REGRESSION ANALYSIS

Regression analysis is a statistical process used to estimate relationships between variables. First, we will use examples involving a limited number of variables to illustrate the concepts of linear regression. Statisticians typically work with a small number of variables, such as demographic information on a population of citizens. In machine learning settings, the number of variables may be much larger, but the basic principles still apply.



Extra resource

Please consult the highly recommended additional readings on **Discrete Maths** and **Calculus** to develop a strong understanding of the mathematical concepts underlying linear regression.

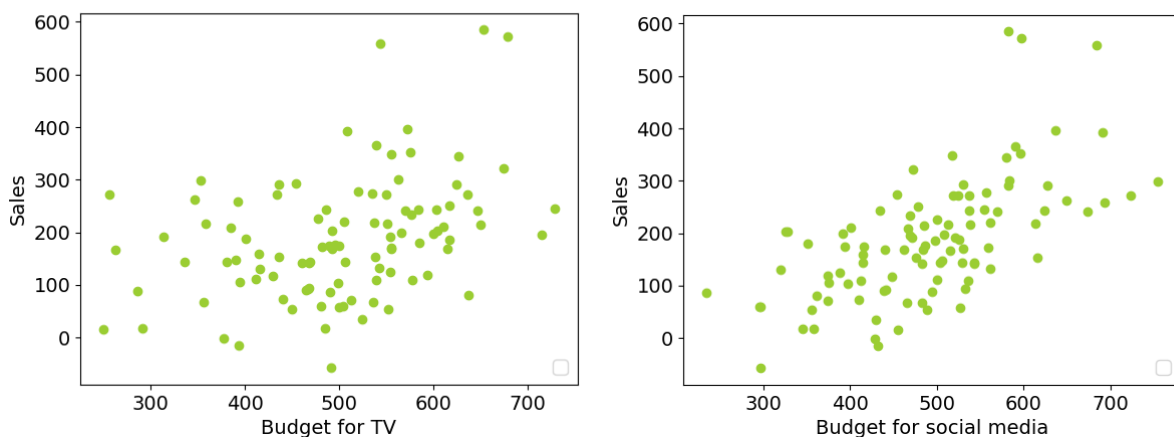
These subjects provide important foundational knowledge, including linear algebra and derivatives, which will enhance your overall understanding of machine learning.

SIMPLE LINEAR REGRESSION

Suppose that we have been asked by a client to give advice on how to advertise their product most effectively. The client offers us some data on which to base our recommendation. They have the sales and advertising budgets of the product in 200 markets, where the advertising budget is split into television ads and ads on social media.

Common sense suggests that spending more money on advertising will increase sales and that different kinds of advertising do so at various rates. Indeed, as seen

in the graphs below, the data show that the higher the budget, the higher the sales.



However, it is hard to tell what the difference in impact between TV and social media ads is. Simple linear regression lets us quantify this difference.

We start by expressing our assumption of a relationship between sales and advertising in the following general mathematical form:

$$Y = f(X)$$

Here Y represents sales, and X is the marketing budget, divided into TV and social media budgets (X_1, X_2). The unknown function f takes these variables and performs mathematical operations that convert the values for X into the values for Y .

Simple linear regression proposes the following specification of f , which approximates the equation of a straight line:

$$Y \approx \beta_0 + \beta_1 X_i$$

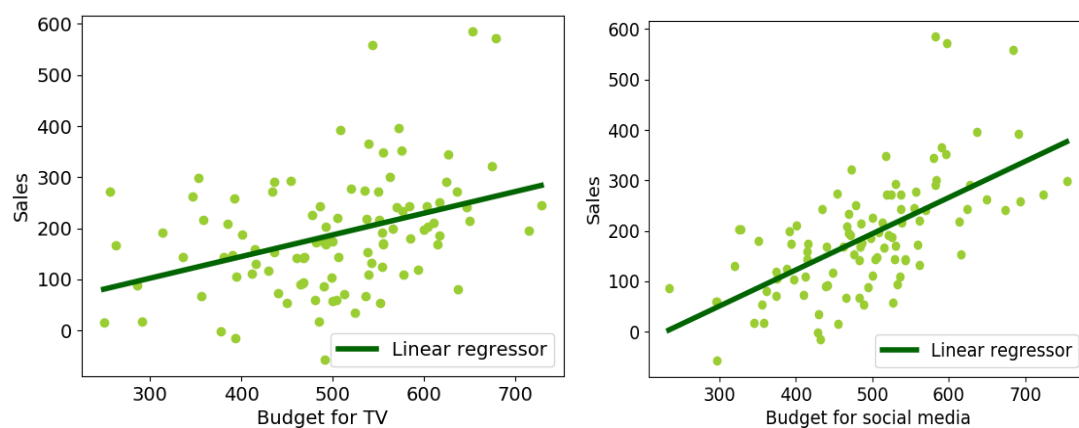
Regardless of which value you specify for β_0 (the intercept) and β_1 (the slope), values produced by this equation will fall along a straight line. The intercept models how many units of the product are sold when there is no money assigned to advertising. The coefficient models how much the sales increase per each single-unit increase in the advertising budget.

The purpose of simple linear regression is to find the straight line that “best fits” the data. The best fit here refers to values for β_i that leave a minimal difference between the straight line produced by f and the observed data. There exist

formulae that determine at what intercept and slope this difference has been minimised.

A minimised difference is *still* a difference: after linear regression, there is still some difference between observed values for \mathbf{Y} , and values for \mathbf{Y} predicted by \mathbf{f} . If your data, when plotted, does not seem to fall along a straight line, linear regression is not the right model for the problem. But even if it does, the straight line is only a model of the data, hence the use of the symbol “ \approx ”.

Lines fitted to our toy data look as follows:



These lines tell us that, according to the data at our disposal, sales are increased by social media advertising more than by TV advertising.

Note that there are fewer instances in our data at the higher end of the x-axis. The assumption that sales will increase linearly may not hold beyond this point if we were to continue to increase the x-value. It is, however, a reasonable approximation for the range of values we have. Moreover, this tried-and-true algorithm is easy to understand, easy to perform, and can provide valuable information.

Let's take a look at a very simple example

The diabetes data set from scikit-learn (**sklearn**). This data set includes two different features:

- **s1**: total serum cholesterol
- **s2**: low-density lipoproteins (a type of cholesterol)

By the nature of these features, there is naturally a correlation between them. This makes for a nice bit of practice in using linear regression.

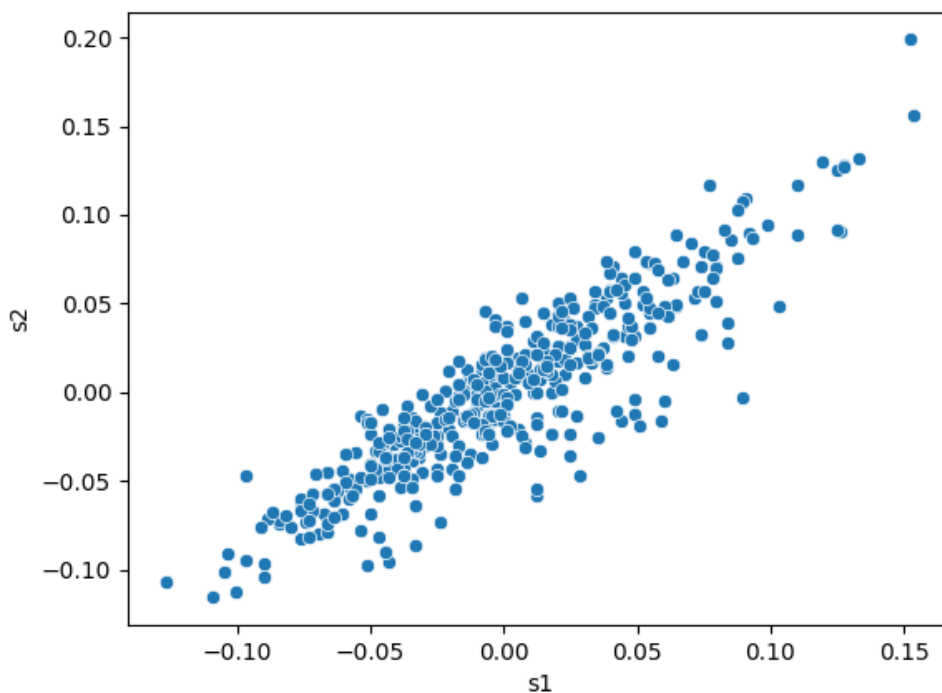
This can simply be loaded as a Pandas dataframe.

```
from sklearn.datasets import load_diabetes
import seaborn as sns
import matplotlib.pyplot as plt

# Load data set
df = load_diabetes(as_frame=True).data[['s1', 's2']]

# Plot data
plt.figure()
sns.scatterplot(data=df, x='s1', y='s2')
plt.show()
plt.close()
```

This code above gives the following scatterplot of the data:



To create our model, we import **LinearRegression** from the **linear_model** module in **sklearn**:

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

To allow compatibility, we must ensure that our model inputs are in the shape **(n_samples, n_features)**.

Because there is only one feature, it will be `(n_samples, 1)`:

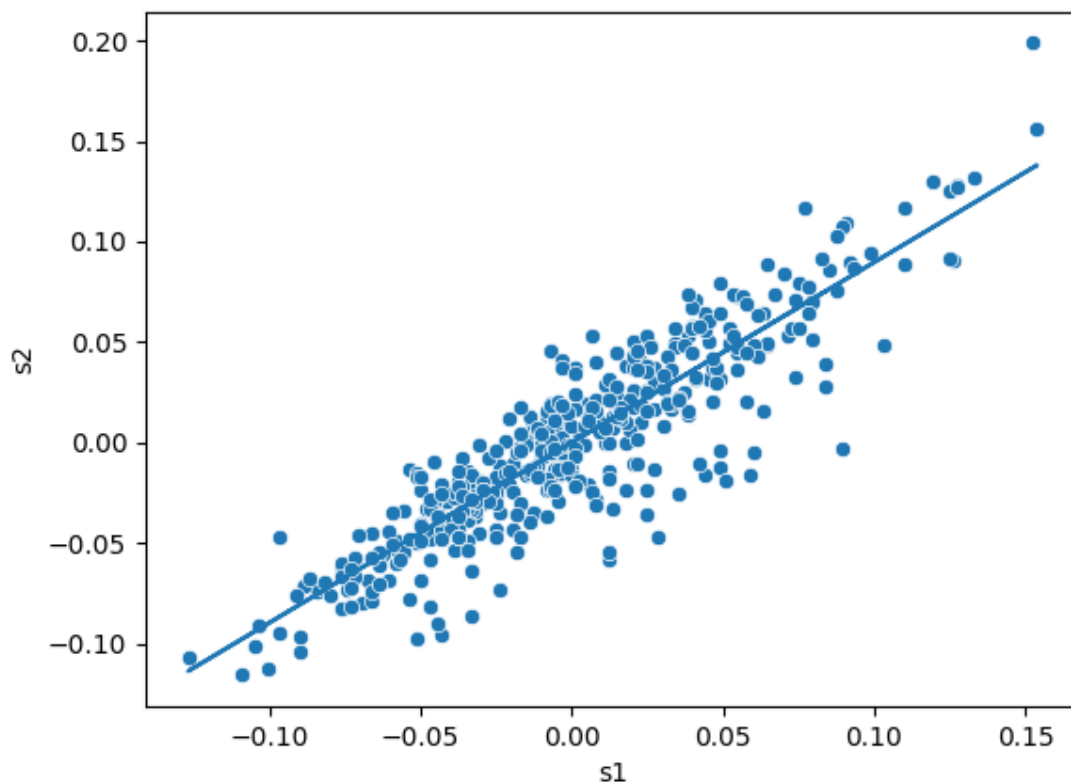
```
X = df['s1'].values.reshape(-1, 1) # create a 2D array
y = df['s2'].values
```

Now that we have our data in the desired shape, we can fit the data and make a prediction:

```
# Fit and predict
simple_model = LinearRegression()
simple_model.fit(X, y)
y_pred = simple_model.predict(X)
```

Let's take a look at our model performance using a plot:

```
# Plot model and data
plt.figure()
sns.scatterplot(data=df, x='s1', y='s2')
plt.plot(X, y_pred)
plt.show()
plt.close()
```



You can see a clear line of best fit going through the data. This was done with just a few lines of code!



A note from the
HyperionDev Team

The difference between machine learning and other statistical and mathematical approaches, such as data mining, is another popular subject of debate. Put simply, while machine learning uses many of the same algorithms and techniques as data mining, one difference lies in what the two disciplines predict.

Data mining discovered previously unknown patterns of knowledge, whereas machine learning is used to reproduce known patterns and knowledge, automatically apply that to other data, and then automatically apply those results to decision making and actions.

Data mining efforts have become extremely prevalent among gamers. Data miners continuously inspect logic in new version releases of games and/or apps to reveal potential future features to other players. This has allowed players to strategise and structure their game so that they can potentially benefit from new features when they arrive.

MULTIPLE LINEAR REGRESSION

We have explored the relationship between one explanatory variable and the response or dependent variable. In many cases, you will be interested in more than a single input variable. In the example case, there were two independent variables: TV and social media. Multiple Linear Regression allows us to model the impact of both variables on the outcome variable at the same time.

Modelling multiple variables is quite useful. A simple linear regression approach may over- or underestimate the impact of a variable on the outcome because it does not have any information about the impact of other variables. Imagine, for example, that our social media budget and billboards budget were increased simultaneously and sales went up shortly thereafter. Our simple regression models would attribute the entire increase in sales to the single variable they were modelling, which would be incorrect. A multiple linear regression model could make a less biased prediction of how much each type of advertisement contributes to sales.

The extension of simple linear regression is fairly straightforward. Instead of a function for \mathbf{Y} with only one coefficient, the function has coefficients for each variable. In the case of two variables, the formula takes the form:

$$\mathbf{Y} = \beta_0 + \beta_1\mathbf{X}_1 + \beta_2\mathbf{X}_2$$

Note that the coefficients (β_i) will not just be the same values as independent simple linear regression models would return. This extended model will adjust each value according to the relative contribution of each variable.



Take note:

Over time, several myths have emerged surrounding Machine Learning, such as “Machine Learning is just summarising data” or “Learning algorithms just discover correlations between events”. Probably one of the most prominent is that “Machine Learning can’t predict unseen events”. This line of thought goes thus: if something has never happened before, its predicted probability must be zero. Right?

On the contrary, machine learning is the art of predicting rare events with high accuracy. If A is one of the causes of B, and B is one of the causes of C, A can lead to C, even if we’ve never seen it happen before. A practical example of the predictive powers of machine learning is the spam filter in your email inbox. Every day, spam filters correctly flag freshly concocted spam emails, based on emails you marked as spam earlier and other predictor data.

Both simple and multiple linear regression can be performed very simply using sklearn. Provided that your data set is properly pre-processed, sklearn infers on its own whether your input has one or multiple features. Fitting looks like this:

```
# Fit a multiple regression model
multiple_model = LinearRegression()
multiple_model.fit(X, y)
```

Let’s apply this to some data on the impact of TV, radio, and newspaper advertising on sales. This data set, **Advertising.csv**, is also in your lesson folder if you want to try it for yourself.

Multiple linear regression applied to this data set returns the coefficients **0.045**, **0.189**, and **-0.001**.

To see what these coefficients say about the variables, it helps to see them in the context of the formula for **Y**:

$$\text{sales} = 2.94 + 0.045(\text{TV}) + 0.189(\text{radio}) - 0.001(\text{newspaper})$$

One thing this shows is that TV and radio advertising has a positive impact on sales, but newspaper advertising has an impact that is close to zero, and even a bit below it. Negative coefficients mean that the variables have an opposite relationship: as one increases, the other one decreases. When comparing the coefficients for TV and radio, we see that the coefficient for radio is larger than that for TV. This means that radio advertising contributes more to total sales than TV advertising does.

Based on these findings, we may recommend a focus on radio advertising. We might also recommend that the newspaper advertising budget be scrapped.

TRAINING AND TEST DATA IN MACHINE LEARNING

When you are teaching a student arithmetic summation, you want to start by teaching them the pattern of summation, and then test whether they can apply that pattern. You will show them that $1+1=2$, $4+5=9$, $2+6=8$, and so forth. If the student memorises all the examples, they could answer the question 'what is $2+6$?' without understanding summation. To test whether they have recovered not just the facts but the pattern behind the facts, you need to test them on numbers that you have not exposed them to directly. For example, you might ask them 'what is $4+9$?'

This intuition forms the motivation behind training and testing data in machine learning. We create a model (e.g. regression) based on some data that we have, but we do not use all of our data: we hide some data samples from the model and then test our model on the hidden data samples to confirm that the model is making valid predictions as opposed to reiterating what was given to it in the training data set.

A **training set** is the actual data set that we use to train the model. The model sees and learns from this data. A **test set** is a set of withheld examples used only to assess the performance of a model. Although the best ratio depends on how much data is available, a common split is a ratio of 80:20. For example, 8000 training examples and 2000 test cases. Sklearn allows us to do this quite easily:

```
# Import the train_test_split functionality from sklearn
from sklearn.model_selection import train_test_split

# Pass your x and y variables in the function and get the train and test
# sets for the X and y variables (4 variables).
# The test_size parameter determines how the data is split, 0.25 means 25%
# of the data will be in the test set.

X_train, X_test, y_train, y_test = train_test_split(X, y, test_size = 0.25)
```

One thing to watch out for when dividing the data is that the test and training set should not be systematically different. If the total data set is ordered in some way, for example by alphabet or by the time of collection, the test set might contain different kinds of instances from the training set. If that occurs, the model's performance will be poor, not because it didn't learn from the data effectively, but because it is being tested on a task that differs from the one it was trained for. For this purpose, it is customary to investigate the distribution of labels in your data and make sure they are similar across your training and test data. In the example below, notice that there are no **0** samples in the test set.

```
# Data sample
X = [1,2,3,4,2,6,7,8,6,7]
y = [0,0,0,0,0,1,1,1,1,1]

# Split data
X_train, X_test, y_train,y_test = train_test_split(X, y, test_size = 0.2,
shuffle=False)

# Compare train and test sets for y
print("y_train {}".format(y_train))
print("y_test {}".format(y_test))
```

Output:

```
>> y_train [0, 0, 0, 0, 0, 1, 1, 1]
>> y_test [1, 1]
```

Notice there are no **0** labels in the test set (**y_test**). This can largely be solved by using `shuffle=True` as a parameter in `train_test_split`. This means that the labels are distributed randomly between the training and test sets, so it is less likely that the training and testing data would be systematically different.

Output:

```
>>y_train [1, 0, 0, 1, 0, 1, 0, 1]
>>y_test  [0, 1]
```

Notice there is now a **0** label in the test set. However, there is still the possibility that most samples of one type end up in the training set, without a representative proportion in the test set. This could also result in lower performance than expected. Consider the case where we set **test_set=0.4**, and we get the following:

Output:

```
>>y_train [0, 1, 1, 1, 0, 1]
>>y_test  [0, 0, 0, 1]
```

Notice the labels are not distributed proportionally between **y_test** and **y_train**. Another way to ensure that data is represented equally in both the training and testing data is to make use of the stratify parameter. This ensures that labels are represented in as close to the same proportions as possible in both the training and testing data.

```
# Split data
```

```
X_train, X_test, y_train, y_test = train_test_split(X, y, test_size= 0.5,
shuffle=True, stratify=y)
```

Output:

```
>>y_train [0, 0, 1, 1, 1, 0]
>>y_test  [1, 0, 0, 1]
```

Note that the labels are now distributed in the same proportion across the training and test sets.



Take note:

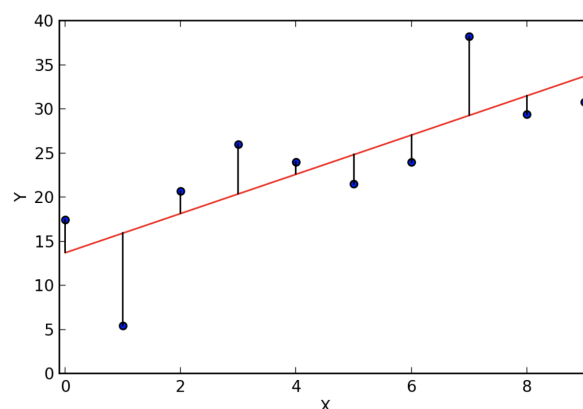
You need to be using **sklearn** version 0.17 or higher to make use of the stratify parameter.

TRAINING AND TEST ERROR

We've discussed before how a regression model is only an approximation of the data. The model predicts values that are close to the observed training values, in hopes of making good predictions on unseen data. The difference between the actual values and the predictions is called the *error*.

The training data set error and the test data set error can be obtained by comparing the predictions to the known, true labels (the gold standard). The test error value is more important as it is the more reliable assessment of the value of the model. After all, we want to use our model on unknown data points, as opposed to applying it to cases for which we already have the actual outcome. In most cases, the training error value will also be lower than the test error value.

There are many different ways to measure the error. As mentioned, the error is the difference between observed and predicted values, as in this plot:



Here we have observed data (**Y**) in dark blue and prediction of a regression model (**y_pred**) along the red line. The error is depicted by vertical black lines. Recall there are a number of different ways one can aggregate the error to get a final score for the model. Two common ones are the root mean squared error (RMSE) and R squared (R^2).



Extra resource

Read more about the [**R-squared metric**](#) to learn how to interpret this metric and the limitations you should be aware of when using it to evaluate regression models.

Instructions

First, read **simple_linear_regression.ipynb** and **multiple_linear_regression.ipynb** to better understand regression analysis and see an example of linear regression in Python.

Compulsory Task 1

Follow these steps:

- Create a Jupyter notebook called **insurance_regression.ipynb**.
- Import **insurance.csv** into your notebook ([Source](#)).
- Use the data in the relevant columns to determine how age affects insurance costs:
 - Plot a scatter plot with age on the x-axis and charges on the y-axis.
 - Using **linear_model.LinearRegression()** from sklearn, fit a model to your data, and make predictions on the data.
 - Plot another scatter plot with the best-fit line.

Compulsory Task 2

Follow these steps:

- Create a Jupyter notebook called **diabetes_regression.ipynb**.
- Read **diabetes.csv** into your Jupyter notebook.
- The **diabetes.csv** aims to predict a person's progression in the condition with respect to various attributes about them.
- Differentiate between the independent variables and the dependent variable and assign them to variables x and y.
- Generate training and test sets comprising 80% and 20% of the data respectively.

- Use a **MinMaxScaler** and **StandardScaler** from **sklearn.preprocessing**. Fit these scalers on the **train set** and use these fit scalers to transform the **train** and **test sets**.
- Generate a multiple linear regression model using the **training set**. Use all of the independent variables.
- Print out the intercept and coefficients of the trained model.
- Generate predictions for the **test set**.
- Compute R-squared for your model on the **test set**. You can use **r2_score** from **sklearn.metrics** to obtain this score.
- Ensure your notebook includes comments about what your code is accomplishing and notes about model outputs such as R-squared.



Rate us
Share your thoughts

HyperionDev strives to provide internationally-excellent course content that helps you achieve your learning outcomes.

Think that the content of this task, or this course as a whole, can be improved or think we've done a good job?

[Click here](#) to share your thoughts anonymously.

REFERENCE(S)

IEEE. (1993). *IEEE Standards Collection: Software Engineering*. IEEE Standard 610.12-1990.