



TASK

Supervised Learning I

Visit our website

Introduction

WELCOME TO THE FIRST SUPERVISED LEARNING 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.



Get in touch
Connect for support

Remember that with our courses, you're not alone! You can contact an expert code reviewer to get support on any aspect of your course.

The best way to get help is to login to Discord at <https://discord.com/invite/hyperdev> where our specialist team is ready to support you.

Our team is happy to offer you support that is tailored to your individual career or education needs. Do not hesitate to ask a question or for additional support!

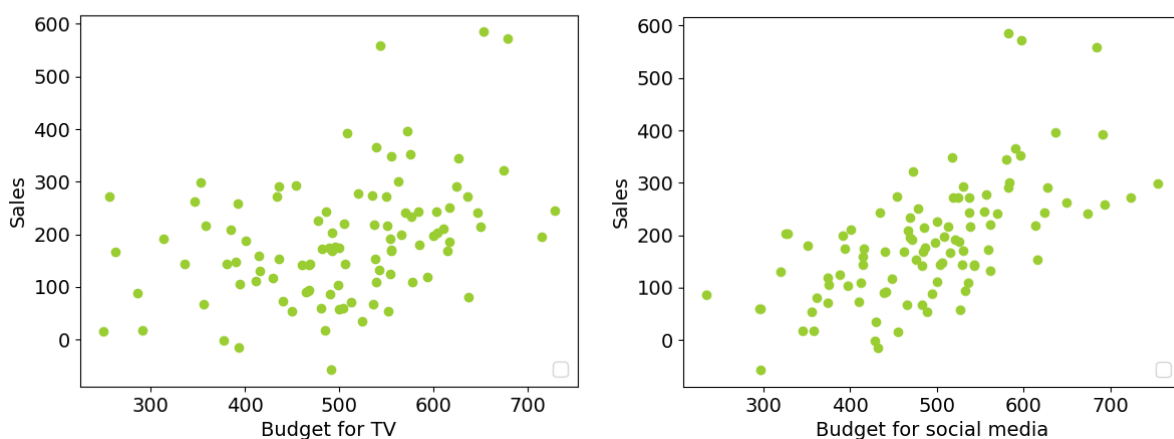
REGRESSION ANALYSIS

Regression analysis is a statistical process used to estimate relationships between variables. There are two types of linear regression: **simple linear regression**, which is the focus of this task, and **multiple linear regression**, which will be taught in the next task. We will teach these concepts using examples involving limited numbers of variables. Statisticians typically do 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.

LINEAR REGRESSION

Suppose that we have been asked by a client to give advice on how to advertise their product most effectively. The client can offer us some data on which to base our recommendation. She offers us the sales and advertising budgets of the product in 200 markets, where the advertising budget is split into the budget for television ads and the budget for 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 are the sales, 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 which 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 Y , and values for Y predicted by 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 ”.



A note from the
HyperionDev Team

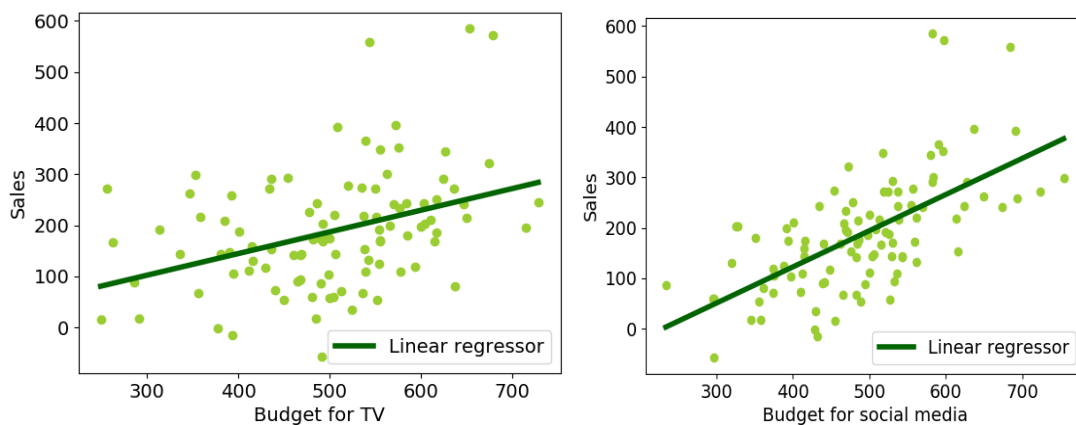
The difference between machine learning and other statistical and mathematical approaches, such as data mining, is another popular subject of debate. Put simply, while

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.

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.

SKLEARN

Scikit-learn (or SKLearn) is a nice and simple starting point for using machine learning. It comes bundled with many techniques that you are already familiar

with (such as preprocessing). However, it also comes bundled with many of the basic machine learning algorithms. It also has a simple and intuitive interface: simply **instantiate** your model, **fit** it to the data and **make a prediction**. It also includes methods for calculating accuracy, loss, etc.!

Let's take a look at a very simple example: the diabetes dataset from SKLearn. This dataset 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.

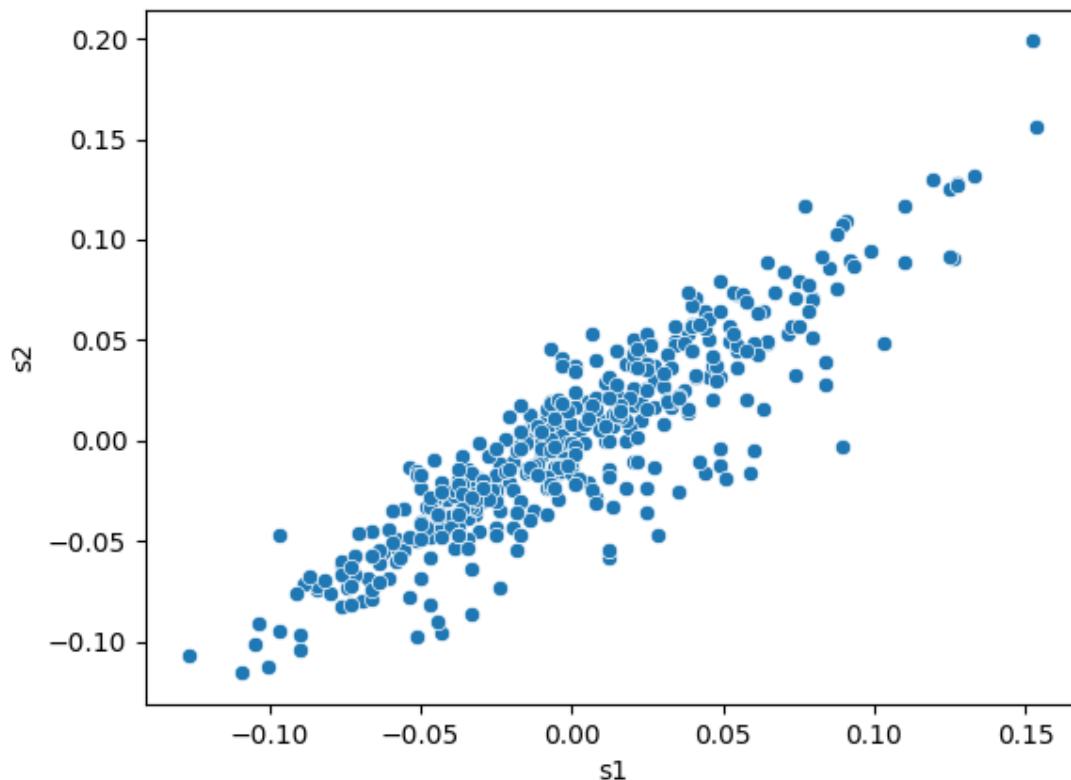
This can simply be loaded as a Pandas dataframe using:

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

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

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

This gives the following scatterplot:



To create our model, we import it from the `linear_model` library in SKLearn:

```
from sklearn.linear_model import LinearRegression

linreg_model = LinearRegression()
```

To allow compatibility, we must ensure that our model inputs are of shape `(n_samples, n_features)`. Because there is only one feature, it will be `(n_samples, 1)`:

```
X = df['s1'].values.reshape(-1, 1)
y = df['s2'].values
```

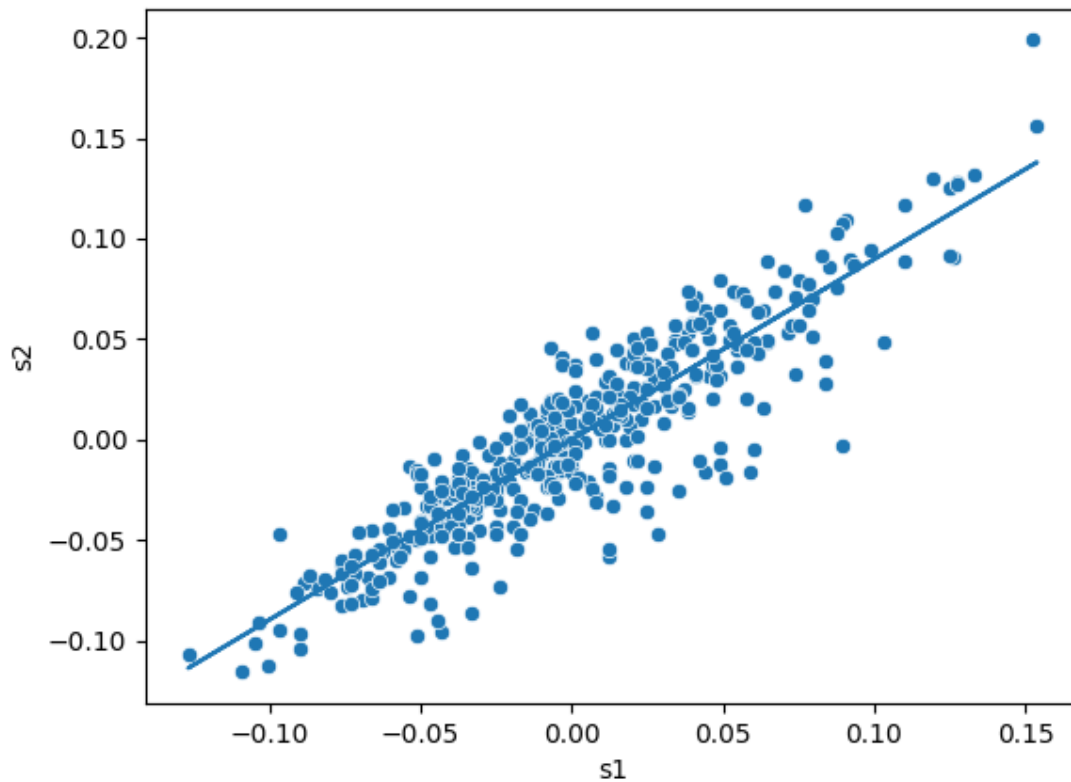
Now that we have our data in the desired shape, we can fit the data and make a prediction:

```
linreg_model.fit(X, y)
y_pred = linreg_model.predict(X)
```

Let's take a look at our model performance:

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

```
plt.plot(X, y_pred)
plt.show()
plt.close()
```



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

Instructions

First, read **[simple-linear-regression.ipynb](#)** to better understand regression analysis and see an example of linear regression in Python.

Compulsory Task 1

Follow these steps:

- 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 data.
 - Plot another scatter plot with the best-fit line.
- Now repeat the above using the data in the relevant columns to determine how BMI affects insurance costs.

If you are having any difficulties, please feel free to contact our specialist team [on Discord](#) for support.

Things to look out for:

1. Make sure that you have installed and set up all programs correctly. You have set up **Dropbox** correctly if you are reading this, but **VS Code** or **Anaconda** may not be installed correctly.
2. If you are not using Windows, please ask a reviewer for alternative instructions.

Completed the task(s)?

Ask an expert to review your work!

[Review work](#)



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.

