

Overview

What is Correlation?

Pearson's Correlation

Spearman's Correlation

Co-variance

What is Correlation?

- Ever wondered how your data variables are linked to one another? Variables within a dataset can be related for lots of reasons.
For example:
 - One variable could cause or depend on the values of another variable.
 - One variable could be lightly associated with another variable.
 - Two variables could depend on a third unknown variable.
 - All of these aspects of correlation and how data variables are dependent or can relate to one another get us thinking about their use. Correlation is very useful in data analysis and modelling to better understand the relationships between variables.
- The statistical relationship between two variables is referred to as their correlation.
- A correlation could be presented in different ways:
 - Positive Correlation: both variables change in the same direction.
 - Neutral Correlation: No relationship in the change of the variables.
 - Negative Correlation: variables change in opposite directions.
- The performance of some algorithms can deteriorate if two or more variables are tightly related, called multicollinearity. An example is linear regression, where one of the offending correlated variables should be removed in order to improve the skill of the model.
- We may also be interested in the correlation between input variables with the output variable in order to provide insight into which variables may or may not be relevant as input for developing a model.
- The structure of the relationship may be known, e.g. it may be linear, or we may have no idea whether a relationship exists between two variables or what structure it may take. Depending on what is known about the relationship and the distribution of the variables, different correlation scores can be calculated.
- In this tutorial guide, we will delve into a correlation score tailored for variables with a Gaussian distribution and a linear relationship. We will also explore another score that does not rely on a specific distribution and captures any monotonic (either increasing or decreasing) relationship.

The Dataset for demo:-

- Before we go diving into correlation methods, let's define a dataset we can use to test the methods.
- We will generate 1,000 samples of two variables with a strong positive correlation. The first variable will be [random numbers](#) drawn from a Gaussian distribution with a mean of 100 and a standard deviation of 20. The second variable will be values from the first variable with Gaussian noise added with a mean of 50 and a standard deviation of 10.
- We will use the `randn()` function to generate random Gaussian values with a mean of 0 and a standard deviation of 1, then multiply the results by our own standard deviation and add the mean to shift the values into the preferred range.
- The pseudorandom number generator is seeded to ensure that we get the same sample of numbers each time the code is run.

Python command is .cov()

An complete example: code is as follows:-

```
# generate data1 and data2
from numpy.random import randn
from numpy.random import seed
from numpy import cov
# seed random number generator
seed(1)
# prepare data
data1 = 20 * randn(1000) + 100
data2 = data1 + (10 * randn(1000) + 50)
# calculate the covariance matrix
covariance = cov(data1, data2)
print(covariance)
```

- Running the example calculates and prints the covariance matrix.
[[385.33297729 389.7545618]
[389.7545618 500.38006058]]
- The covariance between the two variables is 389.75. We can see that it is positive, suggesting the variables change in the same direction as we expect.
- In the example above, since the dataset was contrived (generated manually - with each variable drawn from a Gaussian distribution and the variables linearly correlated) covariance is a reasonable method for describing the relationship.

Pearson's Correlation

- Named after Karl Pearson, The Pearson correlation coefficient can be used to summarize the strength of the linear relationship between two data samples.
- Pearson's correlation coefficient is calculated by dividing the covariance of the two variables by the product of their respective standard deviations.

Formula

>

$$r = \frac{\sum (x_i - \bar{x}) (y_i - \bar{y})}{\sqrt{\sum (x_i - \bar{x})^2 \sum (y_i - \bar{y})^2}}$$

r = correlation coefficient

x_i = values of the x-variable in a sample

\bar{x} = mean of the values of the x-variable

y_i = values of the y-variable in a sample

\bar{y} = mean of the values of the y-variable

The coefficient returns a value between -1 and 1, symbolizing the full spectrum of correlation: from a complete negative correlation to a total positive correlation.

A value of 0 means no correlation.

The value must be interpreted, where often a value below -0.5 or above 0.5 indicates a notable correlation, and values below those values suggests a less notable correlation.

The `pearsonr()` function (from Scipy library) can be used to calculate the Pearson's correlation coefficient between two data samples with the same length.

An example of code - program :-

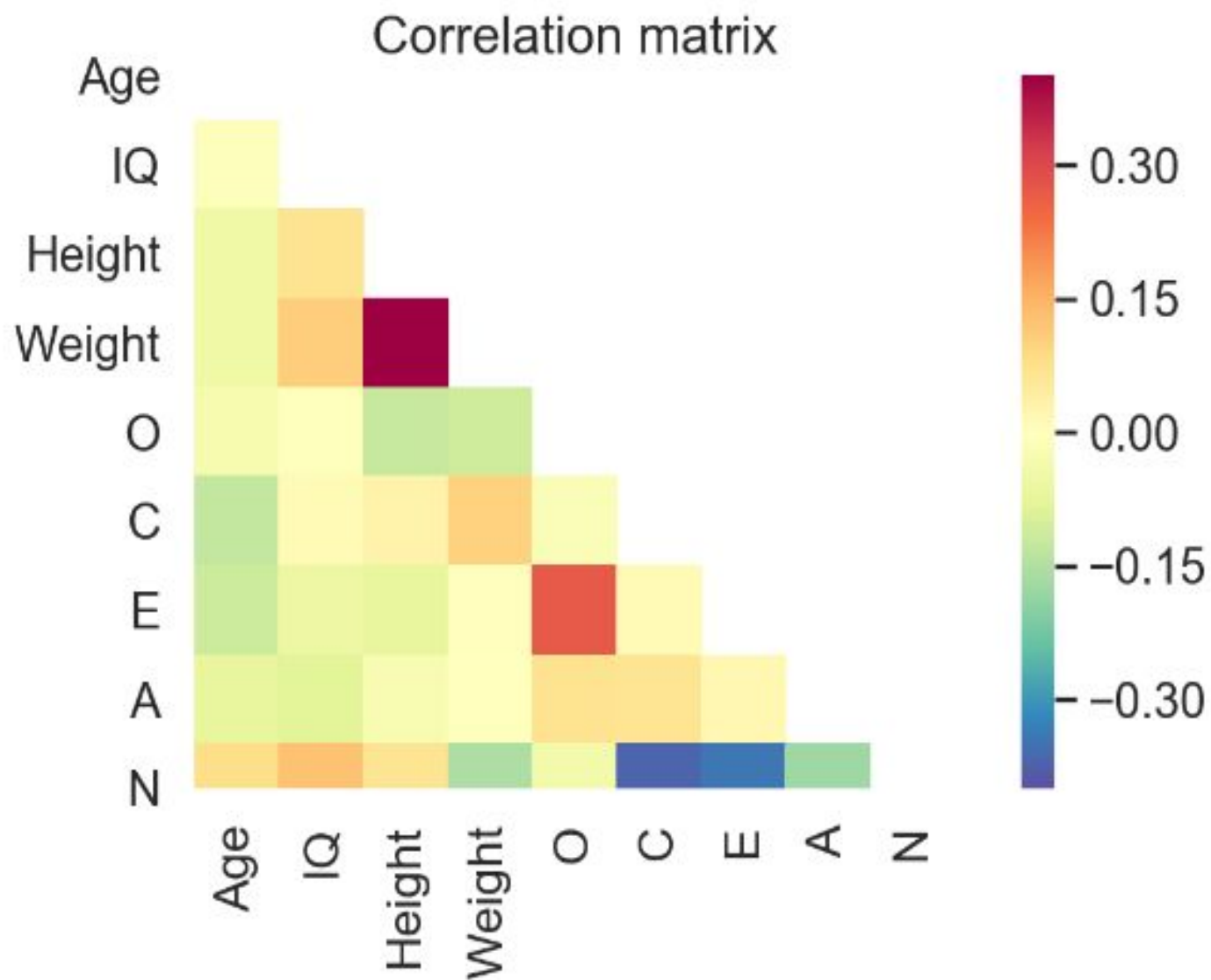
for only 2 columns (features) : data1 & data2,

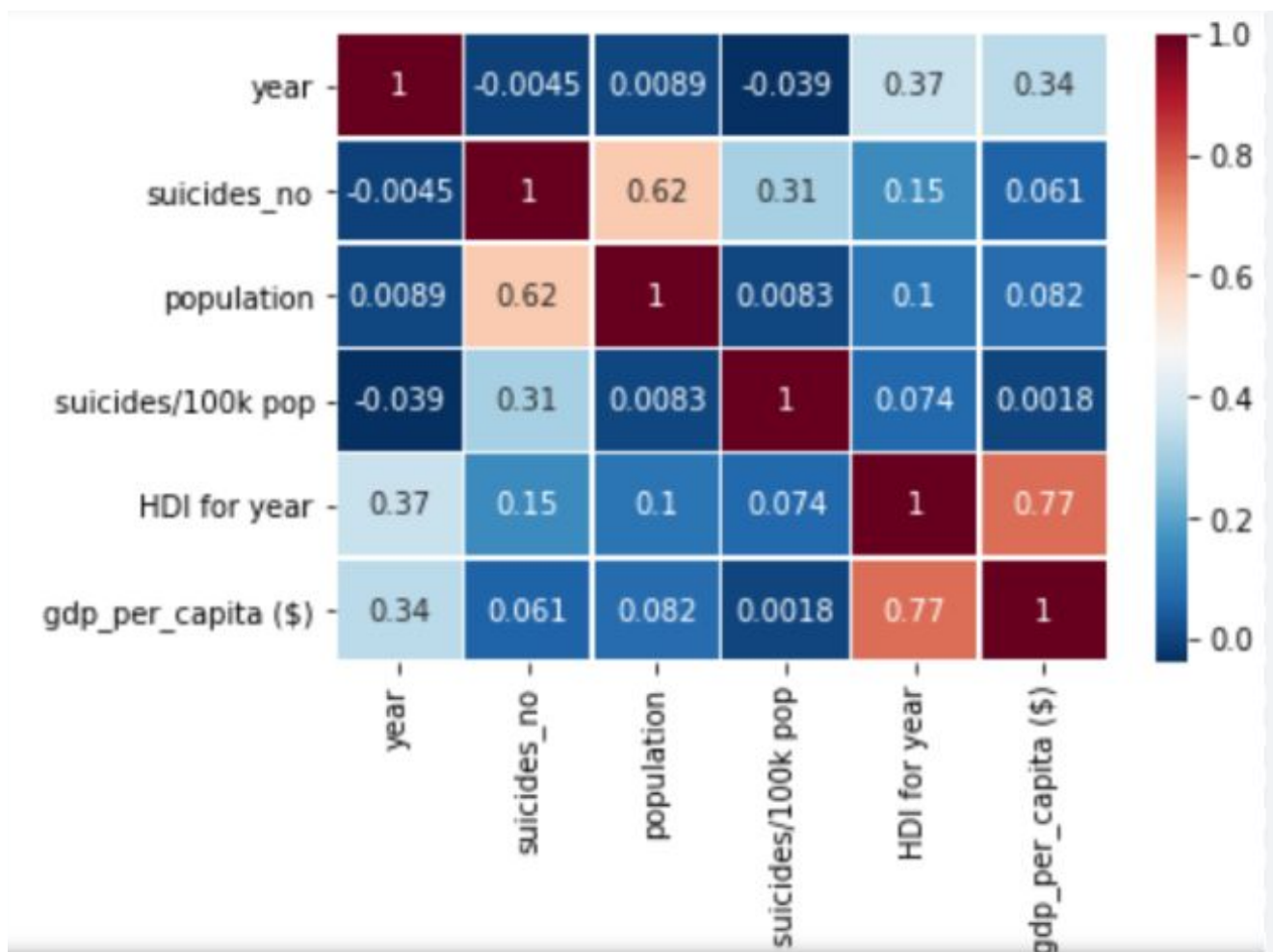
```
from numpy.random import randn
from numpy.random import seed
from scipy.stats import pearsonr
# seed random number generator
seed(1)
# prepare data
data1 = 20 * randn(1000) + 100
data2 = data1 + (10 * randn(1000) + 50)
# calculate Pearson's correlation
corr, _ = pearsonr(data1, data2)
print('Pearsons correlation: %.3f' % corr)
```

Upon execution, Pearson's correlation coefficient is determined and displayed.

The evident positive correlation of 0.888 between the two variables is strong, as it surpasses the 0.5 threshold and approaches 1.0.

Pearsons correlation: 0.888





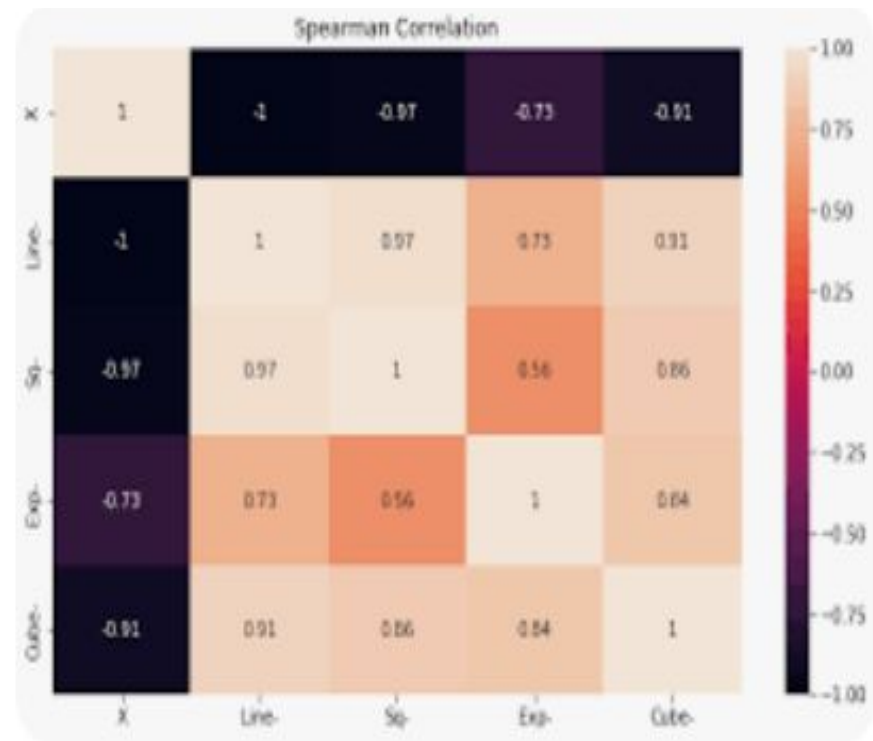
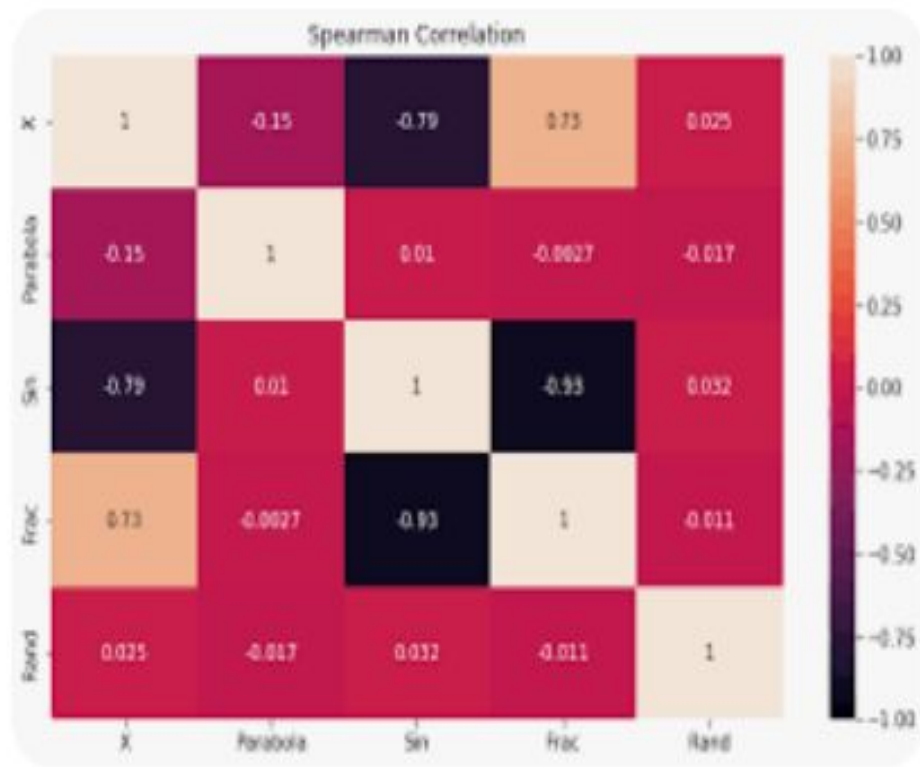


Figure 1: Spearman Correlation

Spearman's Correlation

- Named after Charles Spearman, Spearman's correlation coefficient can be used to summarize the strength between the two data samples. This test of relationship can also be used if there is a linear relationship between the variables but will have slightly less power (e.g. may result in lower coefficient scores).
- **As with the Pearson correlation coefficient, the scores are between -1 and 1 for perfectly negatively correlated variables and perfectly positively correlated respectively.**
- Instead of directly working with the data samples, it operates on the relative ranks of data values.

Formula



$$\rho = 1 - \frac{6 \sum d_i^2}{n(n^2 - 1)}$$

ρ = Spearman's rank correlation coefficient

d_i = difference between the two ranks of each observation

n = number of observations

- If you are unsure of the distribution and possible relationships between two variables, the Spearman correlation coefficient is a good tool to use.

- The `spearmanr()` SciPy function can be used to calculate the Spearman's correlation coefficient between two data samples with the same length.
- We can calculate the correlation between the two variables in our test problem.

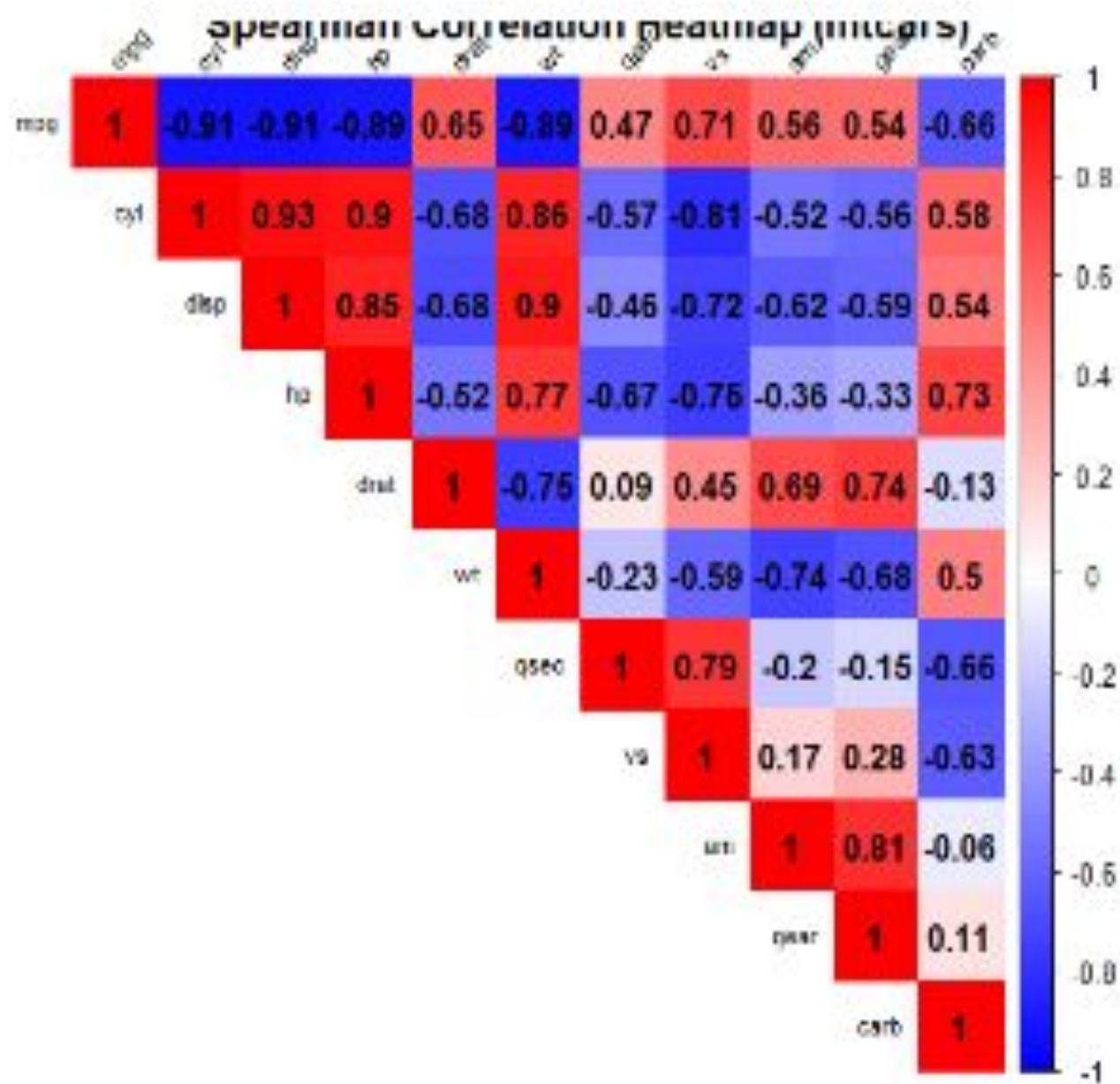
The complete example is listed below.

```
# calculate the spearman's correlation between two variables
from numpy.random import randn
from numpy.random import seed
from scipy.stats import spearmanr
# seed random number generator
seed(1)
# prepare data
data1 = 20 * randn(1000) + 100
data2 = data1 + (10 * randn(1000) + 50)
# calculate spearman's correlation
corr, _ = spearmanr(data1, data2)
print('Spearman's correlation: %.3f' % corr)
```

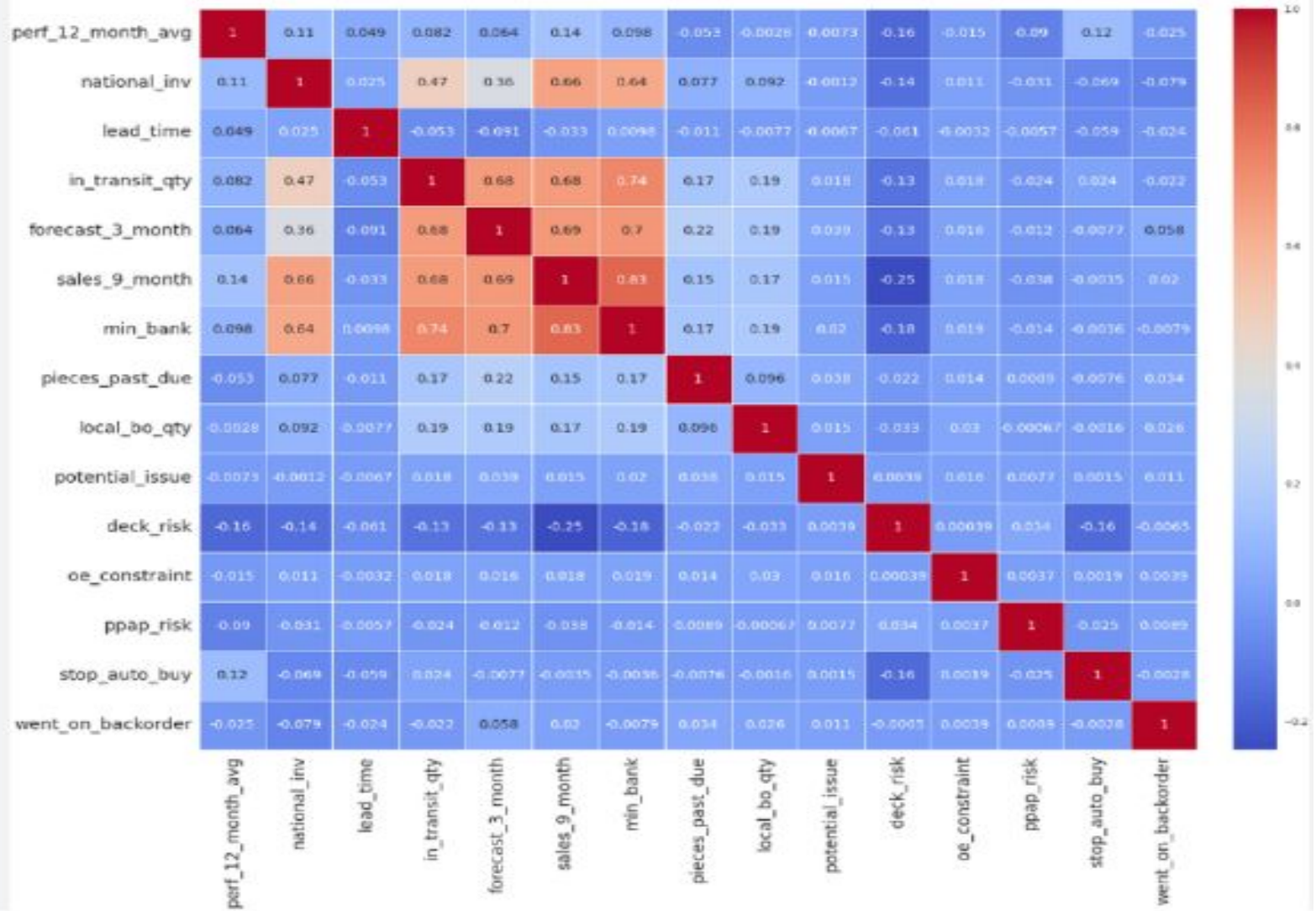
- Running the example calculates and prints the Spearman's correlation coefficient.

Spearman's correlation: 0.872

- As with the Pearson's correlation coefficient, the coefficient can be calculated pair-wise for each variable in a dataset to give a correlation matrix for review.



Spearson's Correlation Heatmap



Covariance

- At the core of many statistical analyses is the concept of linear relationships between variables. This is a relationship that is consistently additive across the two data samples.
- This relationship can be summarized between two variables, called the covariance. It is calculated as the average of the product between the values from each sample, where the values have been centred (by subtracting their respective means).

- The calculation of the sample covariance is as follows:

$$\text{cov}(X, Y) = (\text{sum } (x - \text{mean}(X)) * (y - \text{mean}(Y))) * 1/(n-1)$$

- The use of the mean in the calculation implies that each data should ideally adhere to a Gaussian or Gaussian-like distribution.
- The covariance sign can be interpreted as whether the two variables change in the same direction (positive) or change in different directions (negative). The magnitude of the covariance is not easily interpreted. A covariance value of zero indicates that both variables are completely independent.
- The `cov()` NumPy function can be used to calculate a covariance matrix between two or more variables.

$$\text{covariance} = \text{cov}(\text{data1}, \text{data2})$$

- The diagonal of the matrix contains the covariance between each variable and itself. The other values in the matrix represent the covariance between the two variables; in this case, the remaining two values are the same given that we are calculating the covariance for only two variables.