



Department of Computer Science and Engineering (Data Science)

Sy.B.Tech. Sem: IV

Subject: Statistics For Data Science (DJS23DLPC403)

Experiment 2

Name: Smayan Kulkarni

SAP ID: 60009230142

Date:	Experiment Title: Correlation
Aim	Given a data set of 10 rows. Calculate Karl Pearson's coefficient of correlation, Spearman's rank correlation coefficient (using repeated ranks) manually. Then write a python program to calculate both coefficients and match it with the manually calculated values. Solve the real world problem statements.
Software	Google Colab, Visual Studio Code, Jupyter Notebook
Theory To Be written	What is a strong monotonic relationship? State the types of monotonic relationship with examples.
Implementation	<p>Data:</p> <p>X = 1, 2, 3, 4, 5, 6, 7, 8, 9, 10</p> <p>Y = 5, 6, 7, 8, 7, 9, 10, 10, 11, 12</p> <p>Plot a scatter plot of the above data.</p> <p>Step 2: Karl Pearson's Coefficient of Correlation (r)</p> <p>The formula for Karl Pearson's correlation coefficient is:</p> $r = \frac{n \sum xy - \sum x \sum y}{\sqrt{[n \sum x^2 - (\sum x)^2][n \sum y^2 - (\sum y)^2]}}$ <p>Where:</p> <ul style="list-style-type: none"> • n is the number of data points (in this case 10), • x and y are the individual data points of the variables X and Y.

Step 3: Spearman's Rank Correlation Coefficient (ρ)

Spearman's rank correlation is based on the ranks of the data. The formula is:

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

Where:

- d is the difference between the ranks of the corresponding values of X and Y ,
- n is the number of data points.

Python Code :

1. Write a function to calculate `pearson_correlation(X, Y)`.
2. Write a function to calculate `spearman_rank_correlation(X, Y)`.
3. Use `scipy` to verify Spearman's rank.
4. Print all the three results.

Real world problems.

Q1. The following table gives the data on weekly family consumption expenditure(Y) and weekly family income(X)

$Y :$	70	65	90	95	110	115	120	140	155	150
$X :$	80	100	120	140	160	180	200	220	240	260

(i) Compute the coefficient of correlation between X and Y .

(ii) Test the significance of the coefficient of correlation between X and Y at 5 percent level of significance.

Q2. The following table gives the per capita household expenditure on food (Y) and per capita total household expenditure (X)

$Y :$	60	90	110	125	150	170	180	200	220	230	240	250	255	260	260
$X :$	100	150	200	250	300	350	400	450	500	550	600	650	700	750	800

(i) Compute the coefficient of correlation between X and Y .

(ii) Test the significance of the coefficient of correlation between X and Y at 5 percent level of significance.



	<p>Q3. A company wants to analyze the factors affecting employee productivity. The HR department wants to know:</p> <ol style="list-style-type: none">1. Which independent variable (X1, X2, X3) has the strongest correlation with employee productivity (Y)?2. Is the correlation statistically significant at a 5% level?3. Can we visualize the relationships using scatter plots and a heatmap? <table><tr><th>Employee</th><th>Hours Worked (X1)</th><th>Experience (X2)</th><th>Training Programs (X3)</th><th>Productivity Score (Y)</th></tr><tr><td>1</td><td>35</td><td>2</td><td>1</td><td>50</td></tr><tr><td>2</td><td>40</td><td>3</td><td>2</td><td>55</td></tr><tr><td>3</td><td>45</td><td>5</td><td>3</td><td>65</td></tr><tr><td>4</td><td>50</td><td>7</td><td>2</td><td>70</td></tr><tr><td>5</td><td>52</td><td>9</td><td>3</td><td>78</td></tr><tr><td>6</td><td>55</td><td>10</td><td>4</td><td>85</td></tr><tr><td>7</td><td>60</td><td>12</td><td>4</td><td>88</td></tr><tr><td>8</td><td>62</td><td>14</td><td>5</td><td>90</td></tr><tr><td>9</td><td>65</td><td>15</td><td>6</td><td>92</td></tr><tr><td>10</td><td>68</td><td>18</td><td>6</td><td>94</td></tr><tr><td>11</td><td>70</td><td>20</td><td>7</td><td>96</td></tr><tr><td>12</td><td>75</td><td>22</td><td>8</td><td>98</td></tr></table>	Employee	Hours Worked (X1)	Experience (X2)	Training Programs (X3)	Productivity Score (Y)	1	35	2	1	50	2	40	3	2	55	3	45	5	3	65	4	50	7	2	70	5	52	9	3	78	6	55	10	4	85	7	60	12	4	88	8	62	14	5	90	9	65	15	6	92	10	68	18	6	94	11	70	20	7	96	12	75	22	8	98
Employee	Hours Worked (X1)	Experience (X2)	Training Programs (X3)	Productivity Score (Y)																																																														
1	35	2	1	50																																																														
2	40	3	2	55																																																														
3	45	5	3	65																																																														
4	50	7	2	70																																																														
5	52	9	3	78																																																														
6	55	10	4	85																																																														
7	60	12	4	88																																																														
8	62	14	5	90																																																														
9	65	15	6	92																																																														
10	68	18	6	94																																																														
11	70	20	7	96																																																														
12	75	22	8	98																																																														
Conclusion	<p>Hence, we have calculated Karl Pearson’s coefficient of correlation, Spearman’s rank correlation coefficient (using repeated ranks) manually and have implemented a python program to calculate both. Both the coefficients are matching.</p> <p>Conclusion of real world problem 1</p> <p>Conclusion of real world problem 2</p> <p>Conclusion of real world problem 3</p>																																																																	
Colab Link	https://colab.research.google.com/github/SmayanKulkarni/AI-and-ML-Course/blob/master/SDS/exp_2.ipynb																																																																	

Signature of Faculty