**PRACTICAL-7**

**AIM:** Perform the following apache spark program in DATABRICKS.

| 1. | Find the average number of friends by age. (avgfriends.csv) |
|---|---|
| | <table><tr><th>id</th><th>name</th><th>age</th><th>friends</th></tr><tr><td>0</td><td>Will</td><td>33</td><td>385</td></tr><tr><td>1</td><td>Jean-Luc</td><td>26</td><td>2</td></tr><tr><td>2</td><td>Hugh</td><td>55</td><td>221</td></tr><tr><td>3</td><td>Deanna</td><td>40</td><td>465</td></tr><tr><td>4</td><td>Quark</td><td>68</td><td>21</td></tr><tr><td>5</td><td>Weyoun</td><td>59</td><td>318</td></tr></table> |
| 2. | Use the dataset given and write the code to find the minimum temperature by the location (each whether station) and understand it and modify it to find maximum temperature by the location. (temp.csv) |
| | <table><tr><th>Weather stationID</th><th>Date</th><th>Temp Type</th><th>Temp Value</th></tr><tr><td>ITE00100554</td><td>18000101</td><td>TMAX</td><td>-75</td></tr><tr><td>ITE00100554</td><td>18000101</td><td>TMIN</td><td>-148</td></tr><tr><td>GM000010962</td><td>18000101</td><td>PRCP</td><td>0</td></tr><tr><td>EZE00100082</td><td>18000101</td><td>TMAX</td><td>-86</td></tr><tr><td>EZE00100082</td><td>18000101</td><td>TMIN</td><td>-135</td></tr></table> |
| 3. | Use a given dataset of customers and their spending; find how much amount is spent by the individual customer in total, creating proper RDD in the databricks python notebook and sort out result based on the total spent. (customerorders.csv) |
| | <table><tr><th>CustID</th><th>TrID</th><th>Amount</th></tr><tr><td>44</td><td>8602</td><td>37.19</td></tr><tr><td>35</td><td>5368</td><td>65.89</td></tr><tr><td>2</td><td>3391</td><td>40.64</td></tr><tr><td>47</td><td>6694</td><td>14.98</td></tr><tr><td>29</td><td>680</td><td>13.08</td></tr></table> |
| 4. | Use a text-file given as dataset and count the number of words occur in it. Also, use regular expressions to clean and count the number of words and sort out your output. (wordcount data.txt) |
| |  |