| Team12 | Deliverables |
|---|---|
| **Team12** <ul><li>Title: **Twitter Stream Categorizer**</li><li>Description: The project aims at capturing and analyzing streaming data from Twitter to compare the different categories that come up in tweets. The various categories we want to take into consideration are political, social,entertainment, environmental, science & technology,finance and sports</li></ul><ul><li>Team members:<br>Shriya Sharma        ssharm25<br>Sonal Patil        sspatil4<br>Sai Sameer Tirumalasetti    stiruma</li></ul> | **Deliverables**<br>1. Implement an end to end infrastructure pipeline that takes in twitter data stream into Kinesis, processes it and stores it into S3.<br>2. Perform analysis on the captured data independent of the volume of data being received from Twitter using Kinesis Analytics.<br>3. Giving out near real-time results for the analysis performed updating results of analysis with the changing input data by monitoring the traffic.<br>4. Define classifiers that categorize the data obtained from stream in application and build models based on them.<br>5. Generate graphical visualisation in the model to compare the different categories of the data captured. |
| **Dependencies**<br><br>The following are the dependencies that the mentioned project will be needing to carry out the desired tasks:<br><br>1. Amazon Kinesis Firehose<br>2. Amazon S3 (Storage)<br>3. H20/Amazon ML<br>4. Kibana<br>5. Dataset: Twitter stream<br>6. VCL | **Issues**<br><br>1. Would need a dedicated VCL instance. |