

<p><b>Team12</b></p> <ul style="list-style-type: none"> <li>• Title: <b>Twitter Stream Categorizer</b></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 &amp; technology,finance and sports</li> <li>• Team members: <div> <div>Shriya Sharma</div> <div>ssharm25</div> </div> <div> <div>Sonal Patil</div> <div>sspatil4</div> </div> <div> <div>Sai Sameer Tirumalasetti</div> <div>stiruma</div> </div> </li> </ul>	<p><b>Deliverables</b></p> <ol style="list-style-type: none"> <li>1. Implement an end to end infrastructure pipeline that takes in twitter data stream into Kinesis, processes it and stores it into S3.</li> <li>2. Perform analysis on the captured data independent of the volume of data being received from Twitter using Kinesis Analytics.</li> <li>3. Giving out near real-time results for the analysis performed updating results of analysis with the changing input data by monitoring the traffic.</li> <li>4. Define classifiers that categorize the data obtained from stream in application and build models based on them.</li> <li>5. Generate graphical visualisation in the model to compare the different categories of the data captured.</li> </ol>
<p><b>Status</b></p> <ol style="list-style-type: none"> <li>1. Implemented sample Kinesis stream with 4 shards that is fed with the twitter data of New York area and a consumer that consumes the data.</li> <li>2. Created an S3 bucket that holds the tweet data.</li> <li>3. Implemented Firehose delivery stream that stores the data stream onto S3</li> <li>4. Implemented helper functions to update shards of Kinesis Stream, describe stream, list existing streams, create roles</li> </ol>	<p><b>Issues</b></p> <ol style="list-style-type: none"> <li>1. Handling processing of multiple shards efficiently</li> <li>2. Appropriate partitioning of data</li> <li>3. Determine the optimal number of shards for the kinesis stream.</li> <li>4. Create role policies dynamically through api.</li> <li>5. Distribute shards efficiently through regions when the volume of tweets is scaled up to be received from a larger area.</li> </ol>