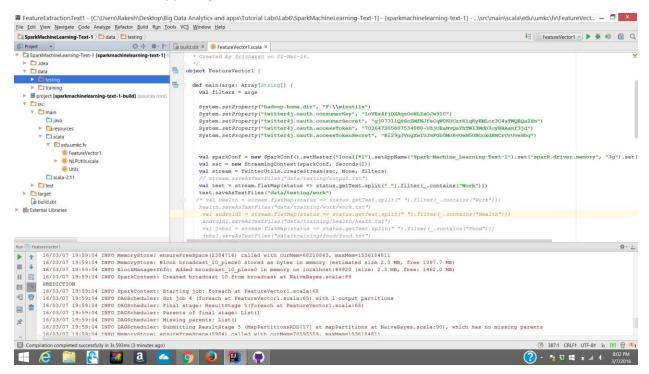# Spark Machine Learning

In this part of assignment large data sets of training and testing data must be collected from Twitter. The content of tweets must be collected in a streaming way.

Here I am collecting data on food, health and work as filters. Divided the data into training and testing data.

Below are the screenshots:

1. PREDICTION is happening.

2. Based on the testing data it is predicting as food.