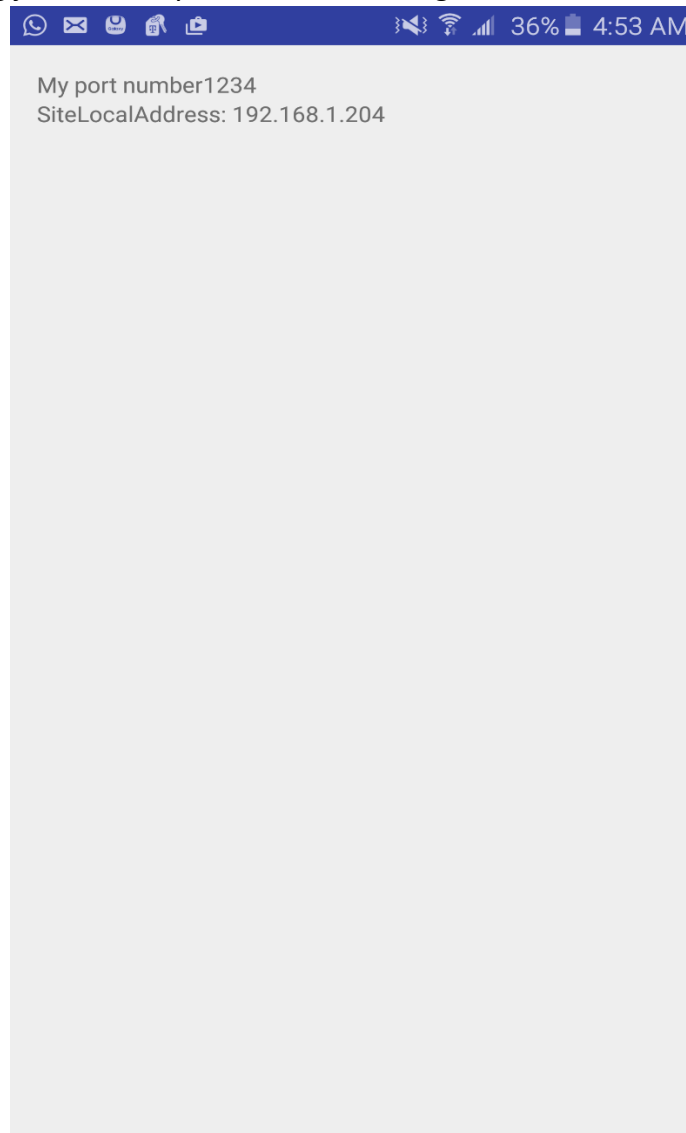# CS5542 Big Data Apps and Analytics
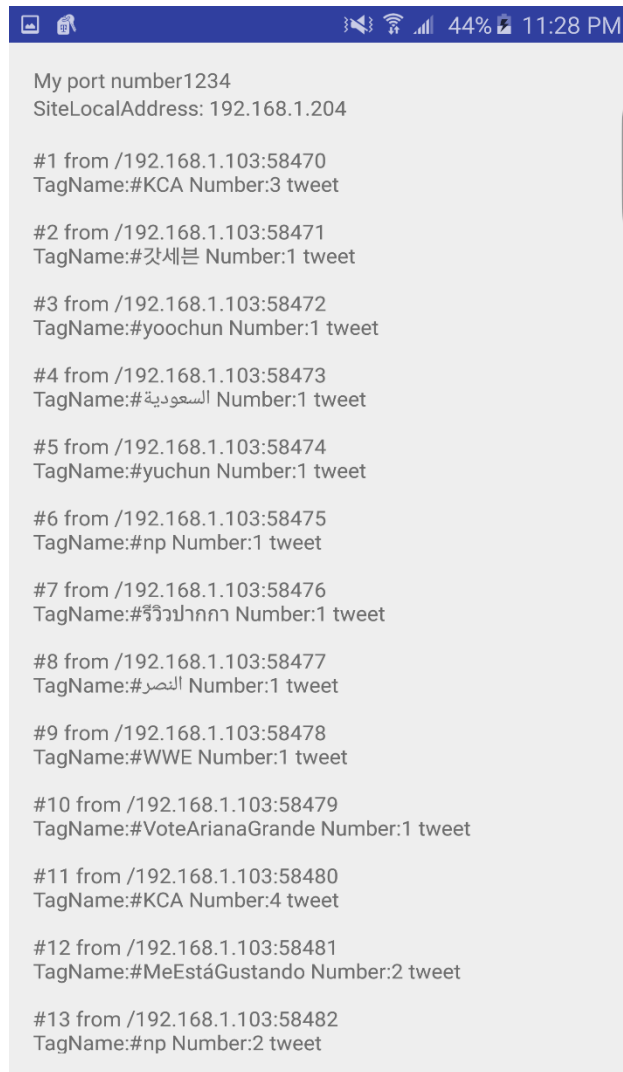# LAB ASSIGNMENT #5 &6
# REPORT and SCREEN SHOTS

**1. Spark and Smartphone/Watch Application Implement a smart application with big data analytics related to your project showing the collaboration between Spark and Smart Apps. Implement Twitter Streaming and perform word count on it and publish the results and showcase it in your Smart Phone/Watch Application.**

1. We open a socket connection between the smart phone and spark . when we run the twitter streaming job, for every 10 seconds hashtags are sent to device.

2. Tweets Top hash tags are displayed on Smart Phone screen

My port number1234
SiteLocalAddress: 192.168.1.204

#1 from /192.168.1.103:58470
TagName:#KCA Number:3 tweet

#2 from /192.168.1.103:58471
TagName:#갓세븐 Number:1 tweet

#3 from /192.168.1.103:58472
TagName:#yoochun Number:1 tweet

#4 from /192.168.1.103:58473
TagName:#السعودية Number:1 tweet

#5 from /192.168.1.103:58474
TagName:#yuchun Number:1 tweet

#6 from /192.168.1.103:58475
TagName:#np Number:1 tweet

#7 from /192.168.1.103:58476
TagName:#รีวิวปากกา Number:1 tweet

#8 from /192.168.1.103:58477
TagName:#النصر Number:1 tweet

#9 from /192.168.1.103:58478
TagName:#WWE Number:1 tweet

#10 from /192.168.1.103:58479
TagName:#VoteArianaGrande Number:1 tweet

#11 from /192.168.1.103:58480
TagName:#KCA Number:4 tweet

#12 from /192.168.1.103:58481
TagName:#MeEstáGustando Number:2 tweet

#13 from /192.168.1.103:58482
TagName:#np Number:2 tweet

**2. Spark ML Lib ApplicationPerform a machine learning algorithm with the Twitter Streaming data tocategorize each Tweet**

**1)Training datasets: Collect different categories of Tweets related to your project.(Categoriescan be based on HashTags /Subjects etc.)**

**2)Test data:theupcoming twitter stream**

**1)** I have collected the different categories of tweets such as Tweets related to restaurants, movies, sports.

I trained my model to predict these categories using the Naïve Bayes Machine learning algorithm

The collected data is Divided into training data and Test data.

Training Data consists of Tweets related to three different categories. I collected the training data using the filters.

Testing data is collection of streaming tweets and we tried to predict the tweets which are more related to which category using classification algorithm.

Training Data and Test Data under Data tab

Tweets in test data are processed using the NLP

Prediction:



Predicting the given test data is more falls under the category of Sports Tweets.