# Increment 2

Group 6: LeeStudents
James Clark
Tirumala Konireddy
Jagadish Tirumalasetty
Wang Zhang
Dr. Yugyung Lee
CS 560: KDM
April 1, 2014

## Data

It turns out there's an ethical approach to consider when crawling sites. Two things to immediately consider: 1) robots.txt usually found at domain.com/robots.txt and 2) the robots meta tag such as <meta name="robots" content="index, follow" />. In the case of this site, it is acceptable to crawl for more data. I will conduct such a crawl at a metered pace, something like one request per 15 seconds. Results should be stored to prevent the need for repeated crawling of the same pages. MongoDB is probably ideal for this. More info is available at
http://www.metatags.info/meta_name_robots
http://www.robotstxt.org/robotstxt.html

## Discovery

- For different algorithms, data format is different[1], therefore we need a data format processor.

Suggested workflow:

Store data set in DB such as HBase or MongoDB.

Get data from DB

Convert data to format

Process with Mahout

- For the classifier to work properly, this set (Naive Bayes classifier) must have at least 50 tweets messages in each category.[2]

The CPSC has an API:
http://www.cpsc.gov/en/Recalls/CPSC-Recalls-Application-Program-Interface-API-Information/
Using this API, I found 5000+ links to recall data in date range 1972-01-01 to 2015-01-01

USA.gov has an API:
http://search.digitalgov.gov/developer/recalls.html

## Algorithms

The CPSC provides an ASMX web service endpoint, and they serve data from 1973 to current. A one-time process is needed to gather all recall data. After that, a periodic check is needed to look for the latest updates.

Connect to CPSC web service, http://www.cpsc.gov/cgibin/CPSCUpcWS/CPSCUpcSvc.asmx
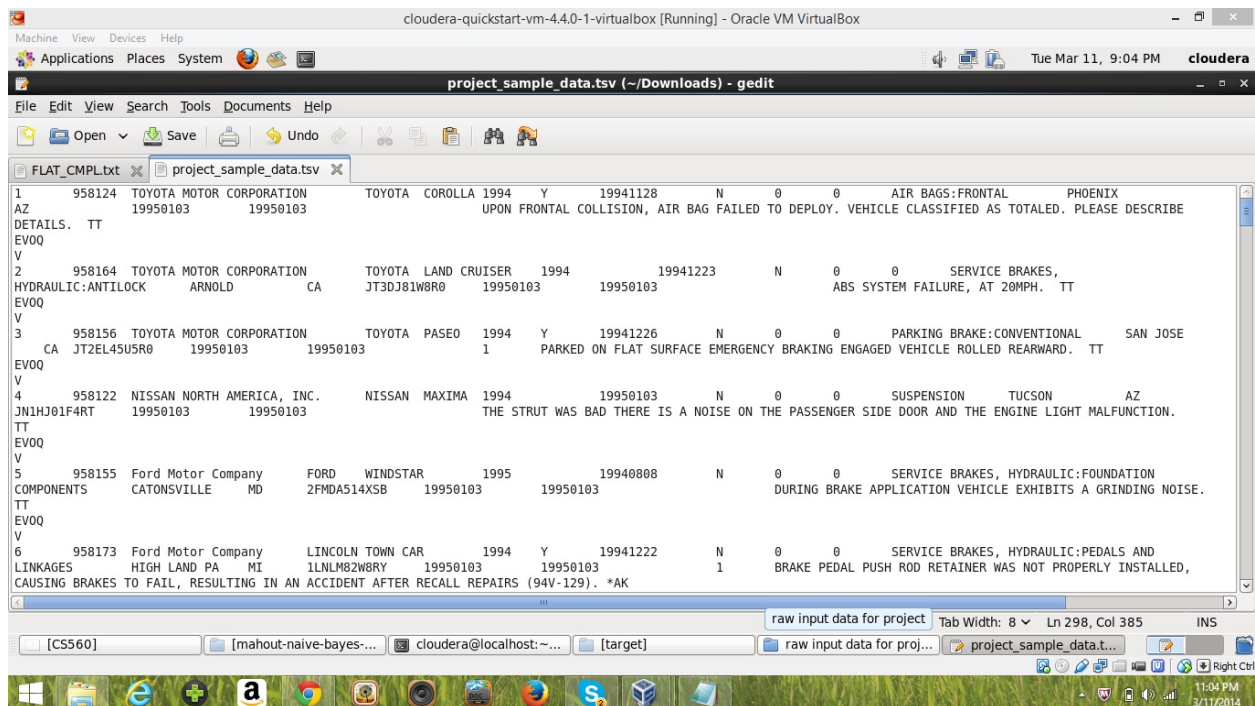Invoke getRecallByDate method
foreach result in list do

---

[1] Feichen, Tutorial 7, Page 18
[2] Feichen, Tutorial 7, Page 19

start json object
store result attributes to json object
visit result url
parse web page at url
store fields from html into json object
return json to caller

Analytical Tasks

1.  Collection of data : We get data from
    http://www-odi.nhtsa.dot.gov/downloads/flatfiles.cfm about all the product recalls that had
    done in past till now. Now, we are manually monitor and download the updated recall data
    flat .txt file.



2. Classification of data :

We are using naive bayes algorithm to classify user comments to a recall
category. We used mahout tool for naive bayes classification. The following steps are
followed to classify data using mahout.
(i) Converting data into sequential files:

Mahout works directly only on sequential files, So to convert this data
which had a tab separator between every two columns in a line, we had developed a java
program "CarsTSVToSeq".

(ii) Upload data :

    After converting data into sequential files, we have to upload the sequential file to hadoop on which mahout directly works on.

(iii) Training and Testing:

    We had developed a Classifier, trainer java files using bayes classes in java to do train the classifier with the given seq files.

We had trained the classifier with varying number of input data records.

We got 90.67% accuracy in testing the trained data entries when we trained 10040 lines of entries, whereas we got only 77.4%, 55.08% accuracy when we trained with 100,000 and 967,487 data records.

(iv) Classification :

    We had used 217 data entries for classification and we found good results by cross-verification.



(v) Uploading results to Solr:

    Because of some issues to work with Glassfish server, we are unable to upload this results to Solr using web restful service.

## Implementation

### Implementation of data model and algorithms (Machine Learning)

We are using flat data tsv files which contains details of about complaint id, manufacturer, model, recalled component description, number of deaths involved in fail of such component, customer description of that component failure etc., in 47 columns using a tab separator between every two columns.

We are using naive bayes algorithm which uses bayes probability by assuming every occurrence of event as independent for training a classifier to classify user comments to a recalled component category. For implementing this naive bayes algorithm we used hadoop files system, mahout tool for classifying data. We had used Solr server to store the classified results so that they will be available for accessing by end user. We are developing a web service to convert the classified results into json format and upload them to Solr server using Glassfish server.

### Application Interface

Most Likely finish the Application interface, furthermore we will change xml to json or HTML format. Here is some screenshots for our app.

choose one of category for further action

Recall Categories

Recent Recall Records

free for personal use

RecallApp_project

Foods

Drugs

Car Components

Supplies

Others

free for personal use

RecallApp_project

**Lao Thai Nam Corp. Voluntarily Recalls Number One Sompa Salted Fish, Because of Possible Health Risk**
date: Mon, 31 Mar 2014 18:37:00 -0400

**Vita Food Products Issues Voluntary Recall of ELF Herring Fillets in Wine Sauce Containing Undeclared Milk**
date: Mon, 31 Mar 2014 10:12:00 -0400

**Fresh Express Issues Recall of Limited Quantity of Already Expired Italian Salad Due to Possible Health Risk, No Illnesses Cited**
date: Fri, 28 Mar 2014 22:45:00 -0400

**Nova Products, Inc. Issues Voluntary Nationwide Recall of Dietary Supplements with Undeclared Active Pharmaceutical Ingredients**
date: Fri, 28 Mar 2014 19:45:00 -0400

**Glaxosmithkline Recalls Alli&reg;**
date: Fri, 28 Mar 2014 12:29:00 -0400

**Macadamia Nut Allergy Alert and Voluntary Recall of 15-Count Boxes of Chocolate Chunk LUNA Bars Due to Package Mislabeling**
date: Thu, 27 Mar 2014 20:49:00 -0400

**BBM Chocolate Distributors, Ltd. Issues Allergy Alert on Undeclared Milk in "CHOCOLAT Alprose 52% CACAO PREMIUM DARK CHOCOLATE, Alprose Swiss Chocolate Dark Chocolate and CHOCOLAT Alprose NAPOLITAINS SWITZERLAND"**
date: Thu, 27 Mar 2014 12:32:00 -0400

**Vermont Common Foods Issues Allergy Alert on Undeclared Peanuts**
date: Wed, 26 Mar 2014 16:55:00 -0400

**Oscar's Smokehouse, Inc. Recalls "Eleven Varieties (11) Of Cheese Spreads" Because of Possible Health Risk**
date: Wed, 26 Mar 2014 15:40:00 -0400

free for personal use

RecallApp_project

**Lao Thai Nam Corp. Voluntarily Recalls Number One Sompa Salted Fish, Because of Possible Health Risk**

Lao Thai Nam Corp., of Dallas, Texas is recalling Number One Sompa Salted Fish, because it has the potential to be contaminated with Clostridium botulinum, a bacterium which can cause life-threatening illness or death. Consumers are warned not to use the product even if it does not look or smell spoiled.

Mon, 31 Mar 2014 18:37:00 -0400

http://www.fda.gov/Safety/Recalls/ucm391226.htm

free for personal use