# Tweeter Analysis on University of Chicago

Presented By

Joan Zhang

THE UNIVERSITY OF CHICAGO

# Outline

- Executive Summary

- Data Source and Exploration

- Methodology & Model Output

- Recommendations

# Executive Summary

- Twitter is a favorite source of text data for analysis.

- The project is to get deeper insights into tweeter followers' opinions towards university of Chicago by using IBM SPSS modeler.

- From a cluster of tweets, we'll find Key differentiators between university of Chicago and other top universities nationwide such as Harvard and Massachusetts Institute of Technology (MIT).

- Get sentiment analysis on University of Chicago and find what distinguishes University of Chicago from other universities such as northwestern and University of Illinois (UIC) in the Chicago metropolitan area.

- Provides recommendation how to boost University of Chicago's name recognition

- Keywords: Tweeter, IBM SPSS modeler, university of chicago(uchicago)

# Data Source and Exploration

## Collect Initial Data & Explore data

- A collection of ~7K tweets containing university of Chicago , and other universities from tweeter

- Tweeter data comes in JSON format. After combining all file, it is roughly 21 M records

- Below chart is an example of tweeter



## Data Preparation

- Below chart represents the variables we selected for this project



- Elephant bird package was used to eliminate the messages not related to the University of Chicago. The final records come to 15,886 in a text file



- Retweet count (RT) is added as a new column to separate the original messages. RT is given



THE UNIVERSITY OF CHICAGO

# Methodology

- IBM SPSS is used to analyze tweeter contents and create two types of models. Basic model [1] and sentimental analysis model

- Basic model with opinion loaded of resource template to compare University of Chicago with top universities nationwide and other universities in metropolitan area

- Model with sentimental analysis package loaded to analyze retweeters' sentiment towards University of Chicago

- Cluster tweets in terms of users' retweet count and users' location

1 please see appendix on basic model chart

THE UNIVERSITY OF
CHICAGO

# Model output on UChicago

THE UNIVERSITY OF CHICAGO

# Key Difference Between University of Chicago and Other Top Universities Nationwide and Other Ones in Chicago Metropolitan area[2]

- UChicago has few tweets in foreign languages and less number of tweets compare with top university nationwide
- Uchicago more related to well-known individual Vs Northwestern and UIC

| University Name | % of top three languages of tweets distribution | | | Total # R C | Highest # of tweets | Descriptors & Documents | | Extraction of top concepts related to excludes itself | | |
|---|---|---|---|---|---|---|---|---|---|---|
| Uchicago | English 94% | Spanish 2.68 | Unknown 0.67% | 6,985 | 1,979 | 197 | 3,541 | Obama | People | http |
| Harvard | English 74.82% | Spanish 8.82% | French 2.72% | 8,157 | 2,241 | 68 | 182 | Evacuates | Build | Word master &leadership |
| MIT | German 78.84% | English 9.33% | Turkish 3.55% | 24,916 | 10,932 | 244 | 537 | Science | Innovation | International-HK |
| NWestern | English 94% | German 1.36% | Romanian 0.7% | 60,000 | 33,631 | 782 | 28157 | Sports | geographical location | Academy |
| UIC | English 77.25% | German 7.02% | Spanish 4.08% | 22,481 | 7,500 | 283 | 10,826 | geographical location | Sports | Sex |

THE UNIVERSITY OF CHICAGO

# Sentiment Analysis towards UChicago

- There are 646 records coming out from tweeter data **POSITIVE**[3] side. We see that people saying happy, in love and express good feeling about Uchicago.

- 7.5% response in the category occupation, 6.2% response in the category academy and 4.7 % in geographical location, 4.3% our 4.7% is from America.

- We see that there are 438 records from tweeter data on **NEGATIVE** side. People express their negative feeling.

- 10.5% response in the category of academy, 7.3% in human resource and 6.2 is in occupation.

THE UNIVERSITY OF CHICAGO

# Recommendation to Boost Uchicago's Name Recognition

- Invite President Obama giving speech as people did not know Obama taught in Uchicago. Below is what people said from tweeter.

| 51 | RT @debbie_mayo: People who say Obama doesn't know the Constitution realiz he literally taught Constitutional law at UChicago right? #GOPDâ€¦ |

- Organize a tour via trauma center. People are excited about newly build trauma center and want to learn. See tweeters data.

RT @constantnatalie: omg omg omg UChicago will open a level 1 trauma center at UChicago Med!!!!!!! @TraumaCenterNow WON! https://t.co/3THWmâ€¦

.@UChicago to build Level 1 trauma center at @UChicagoMed center in Hyde Park – huge news for South Side community https://t.co/OB7IWSxP6x

NBC Chicago: UChicago Medicine Takes First Step Toward New Level 1 Trauma Center https://t.co/NFkbqgAiQD

In past 10 days, 2 NTT union elections won, @uchicagogsu won raises & @TraumaCenterNow wins trauma center at UChicago. Congrats to all!

- Hold network events on popular academic topics through alumni from different universities as larger alumni tweeters are active
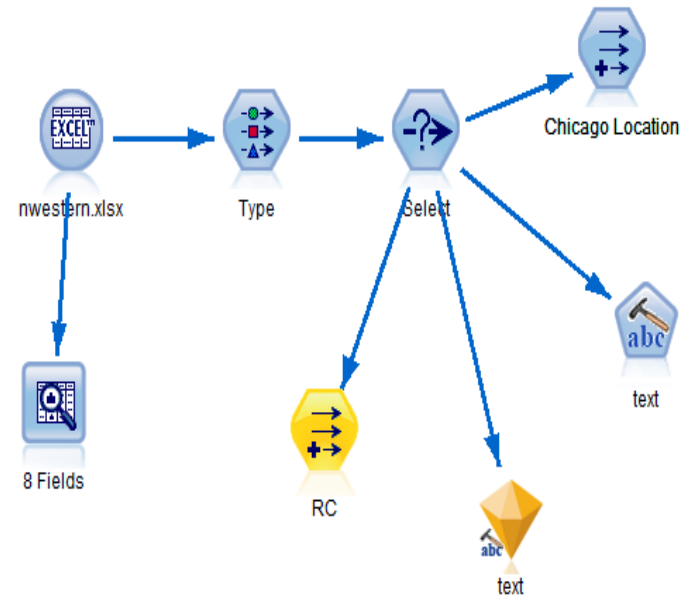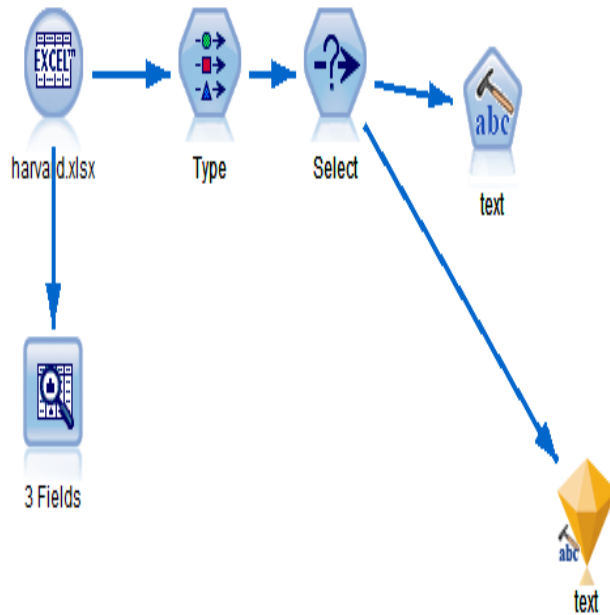
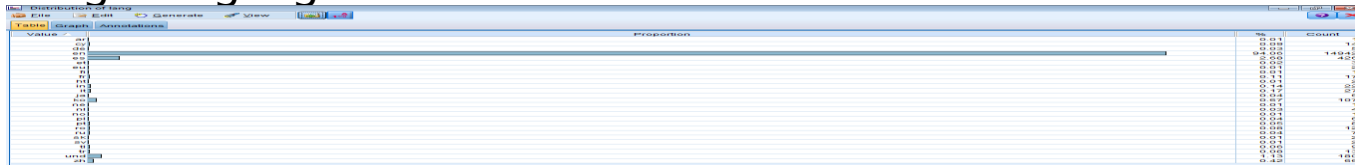| 1 | RT @DUChicagoAlumni: Chicago alums, we're only 18 twitter followers away from 500! Who do you know who'd like @DU_Alumni updates? Tag... htâ€¦ | chicagoan sports/sports by type |

THE UNIVERSITY OF CHICAGO

# Appendix

- Basic Models ( select Harvard and
  Nworthern as examples [1]
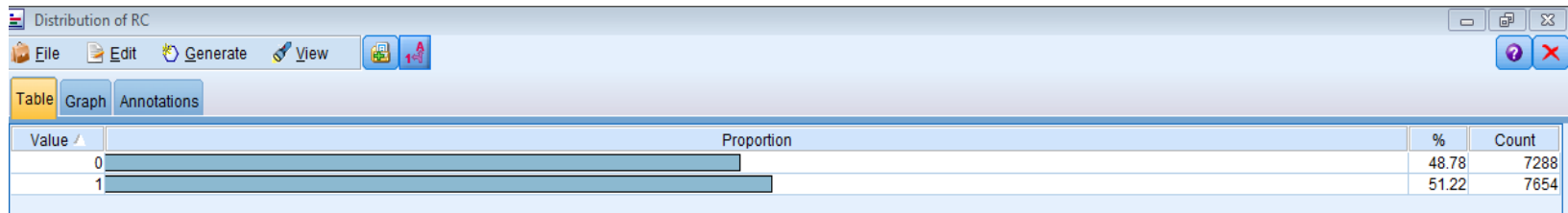
# Chicago detail[2]

- Chicago language distribution chart



- Separate the original messages vs. re-tweets and assign them different weights / counts



- The frequency of retweet count Vs non retweet counts

# Chicago Detail[2]

Concept extracted

Descriptors and Documents in each Categories

# Harvard Detail[2]

- language distribution



- Rewet count is over 8, 000 and right skewed



- Concept extracted



- Descriptors and Categories

# MIT Detail $^2$

- language distribution



- Concept extracted

# NorthWestern$^2$

- language distribution



- Frequency of Rewet count



- Concept extracted



- Descriptors and Categories

# UIC Detail [2]

- language distribution



- Concept extracted



- Descriptors and Categories



THE UNIVERSITY OF CHICAGO

# Frequency of tweets from Metropolitan universities [2]

**Chicago**

Frequency of tweets from 1.38% total 2,185 counts

| Value | Proportion | % | Count |
|---|---|---|---|
| 0 | | 98.62 | 155774 |
| 1 | | 1.38 | 2185 |

**Northwestern**

Frequency of retweet counts is : 51.13% total 80765 counts.

| Value | Proportion | % | Count |
|---|---|---|---|
| 0 | | 48.87 | 77194 |
| 1 | | 51.13 | 80765 |

**UIC**

Frequency of retweet counts is : 46.51% total 25826 counts.

| Value | Proportion | % | Count |
|---|---|---|---|
| 0 | | 53.49 | 29704 |
| 1 | | 46.51 | 25826 |

**Chicago location** : Frequency of tweets: 3.77% total 2094 counts

| Value | Proportion | % | Count |
|---|---|---|---|
| 0 | | 96.23 | 53436 |
| 1 | | 3.77 | 2094 |

THE UNIVERSITY OF CHICAGO

# Clustering and segmentation of users with similar tweets / attitudes toward UChicago

- Highest tweet is 6,622
- Frequency of retweets count is 51.24%

| Value | Proportion | % | Count |
|---|---|---|---|
| 0 | | 46.34 | |
| 1 | | 53.66 | |

# Positive feeling on UChicago [3]

# Added Categories as New Fields



| | | Category_uchicago | Category_uchicago/uchicago and negative | Category_uchicago/uchicago and positive |
|---|---|---|---|---|
| 1 | al harassment in science - my piece on an episode @UChicago. h/t @azeen for leading the way https://t.co/v... | 0 | 0 | 0 |
| 2 | r an Assistant Director of Strategic Events https://t.co/wesfsmLGE1 #JobsMonday @UChicago | 0 | 0 | 0 |
| 3 | $50 million donation to benefit low-income students with high potential https://t.co/d89RMcGDwn https://t.câ€¦ | 0 | 0 | 0 |
| 4 | ne at @UChicago is under enormous pressure to STFU. To save time, reputation, career. Thinking anything ... | 0 | 0 | 0 |
| 5 | uld be fantine and uchicago would be felix tholomyÂ¨s | 1 | 0 | 0 |
| 6 | rsity of Chicago today: @dapizzadude1 @aztec_daves @theroostkitchen | 0 | 0 | 0 |
| 7 | u would follow me back | 0 | 0 | 0 |
| 8 | t Donates $50 Million for University of Chicago's 'Lower-Inco... Read more: https://t.co/SrubxpUcKi | 0 | 0 | 0 |
| 9 | icago &amp; another one of our $2,500 winners. Read more: https://t.co/qEE8EoOBsU https://t.co/uQlauT5A... | 1 | 0 | 0 |
| 10 | walk to the library. #UChicago #Hogwarts https://t.co/EOa1cMmCb4 https://t.co/X4z8wnEQQm | 0 | 0 | 0 |

THE UNIVERSITY OF CHICAGO

# Rules on Chicago sentiment

- Add rules

# Pig script code

- REGISTER hdfs:///jar/elephant-bird/json-simple-1.1.1.jar;
- REGISTER hdfs:///jar/elephant-bird/elephant-bird-core-4.6-SNAPSHOT.jar;
- REGISTER hdfs:///jar/elephant-bird/elephant-bird-pig-4.6-SNAPSHOT.jar;
- REGISTER hdfs:///jar/elephant-bird/elephant-bird-hadoop-compat-4.6-SNAPSHOT.jar;
- REGISTER hdfs:///jar/tutorial.jar;

- A = load '/user/kadochnikov/twitter_full/' USING com.twitter.elephantbird.pig.load.JsonLoader('-nestedLoad') as tweets;
- C = FOREACH A GENERATE
- (CHARARRAY)tweets#'id' as id,
- (CHARARRAY)tweets#'lang' AS lang,
- (CHARARRAY)tweets#'created_at' AS created_at,
- (CHARARRAY)tweets#'text' AS text,
- tweets#'user' as user_info,
- tweets#'retweeted_status' as retweet_info
- ;
- D = FOREACH C GENERATE
- id,
- lang,
- created_at,
- (CHARARRAY)user_info#'screen_name' as screen_name,
- (CHARARRAY)user_info#'name' as name,
- (CHARARRAY)user_info#'location' as location,
- (INT)retweet_info#'retweet_count' as retweet_count,
- REPLACE(text, '\n', ' ') AS text
- ; E = FILTER D BY (text matches '(?i).*uchicago.*'); --(your other filter words go here)

- rmf /user/joanzhang/Test_Tweet

- --forcing a single file output using PARALLEL 1
- reduced = FOREACH (GROUP E BY RANDOM()) GENERATE FLATTEN(E) PARALLEL 1;
- STORE reduced INTO '/user/joanzhang/Test_Tweet' using PigStorage('\t','-schema');

- fs -getmerge -nl /user/joanzhang/Test_Tweet /home/joanzhang/Tweet/reduce.txt

- --on local Linux, run the command to reserver header and remove other staff of pig
- awk 'NR==1 || NR>4 {print}' /home/joanzhang/Tweet/reduce.txt > /home/joanzhang/Tweet/tweet.txt