



[732A47] Text Mining

Text Mining Project

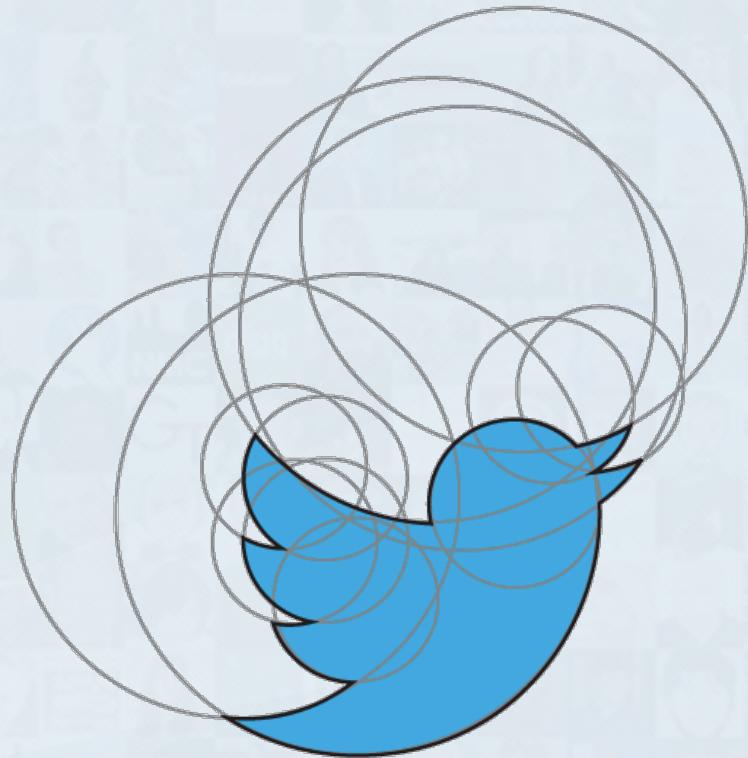
Retweets

student

Sebastiano Milardo

Outline

- *Problem Description and Data*
- *Analysis*
- *Conclusion*



Problem Description

- *What characterizes a tweet that gets retweeted by many people?*

Problem Description

- What characterizes a tweet that gets retweeted by many people?

Successful Example

6 Nov 12

810,000+ Retweets

300,000+ Favorites

Souce: <https://2012.twitter.com/en/golden-tweets.html>

Successful Example

 **Barack Obama** 
@BarackObama

Four more years.
pic.twitter.com/bAJE6Vom

 Risposta  Retweet  Aggiungi ai preferiti  Altro



802.622 RETWEET 300.119 PREFERITI



5:16 AM - 7 Nov 12 Segnala contenuto

Problem Description

Too many...



...too few!

Data

- Google Analytics

- **Queries:** 10-9 Cerrone, Abortion, Acne, Alex Morgan, AllyFollowSpree, Amanda Todd, Apple, Audi, Aurora Shooting, Back Pain, Bank, Barack Obama, Barcelona, Bared to you, Baseball, Basket, Big Bang Theory, Bloody Mary, BMW, Cancer, Cat, Chuck, CNN, Converse, Dad, Dance, Defendig Jacob, Depression, Derrick Rose, Diabetes, Dodge, Dog, Economic, Election, Enviroment, ESPN, Fast & Furious 6, Fifty Shades Darker, Fifty Shades Freed, Flowers, Food, Football, Ford, Fox News, Froch, Galaxy, Game of Thrones, Gangnam Style, Gas, Gay Marriage, Glee, Golf, Good Morning, Heart, Hello, Herpes, Honda, Hotel, Huffington Post, Hunger Games, Hurricane, Hurricane Isaac, Immigration, Instagram, iPhone, Jeremy Lin, Joe Paterno, Justin Bieber, Kate Middleton, Kevin Durant, KJ Noons, Kobe, KONY 2012, Lebron, Lolo Jones, Love, Magic Mike, Marc Gasol, Margarita, McDonald, MentionSomeoneYouLove, Michael Clarcke Duncan, Michael Phelps, Mitt Romney, Mojito, Mom, Money, Morgan Freeman, MSN, Music, NBC, Neil Armstrong, News, Nexus, Nicki Minaj, Nike, Nokia, NY, Obama, Obamacare, ObamacareKidsBooks, Olympics, One Direction, Parking, Paul Ryan, Peyton Manning, Photo, Pizza, Pizza Hut, President, Presidential Election, Prometheus, PSY, Restaurant, Ron Paul, Selena Gomez, Serena Williams, Snooki, Son, Song, SOPA, Sport, Starbucks, Super Bowl, Swimming , Ted, Tennis, ThatOneExWho, The Avengers, The Dark Knight Rises, The Serpent's Shadow, Toyota, Twilight, UEFA Euro 2012 , UFC160, Walmart, WeAllKnowThatOnePersonWho, Whitney Houston, WI, X Factor

Information retrieval

- Topsy.com
- Twitter Streamer API
- Twitter API
 - GET search
 - GET users/lookup
 - GET statuses/user_timeline

Information retrieval

- ~~Topsy.com~~
- ~~Twitter Streamer API~~
- Twitter API
 - GET search
 - GET users/lookup
 - GET statuses/user_timeline

Information retrieval

- 92 010 Unique tweets
 - 9 606 Retweeted
 - 82 404 Non-Retweeted
- Hardware Limits...
 - 9 000 Retweeted
 - 9 000 Non-Retweeted

Features

- Text Mining:
 - Normalization (lowercase)
 - Removing Stop Words
- Features A:
 - Top 600 Most common words
 - Top 300 Most common Screen Names (@)
 - Top 400 Most common Hash Tag (#)
 - Average word length
 - Lexical diversity

Naive Bayes classifier



Results

F-measure: 0.61 - Accuracy: 0.58 - Precision: 0.56 - Recall: 0.66

			F	a	T
			1	s	r
			e	e	e
has(enviroment) = True	False : True =	16.3 : 1.0			
has(prometheus) = True	False : True =	11.1 : 1.0			
has(f1) = True	False : True =	9.1 : 1.0			
has(retweet) = True	True : False =	8.7 : 1.0			
has(sale) = True	False : True =	7.5 : 1.0			
has(mojito) = True	False : True =	7.3 : 1.0			
avg<3 = True	False : True =	7.1 : 1.0			
has(syria) = True	False : True =	6.8 : 1.0			
has(sopa) = True	False : True =	6.5 : 1.0			
has(mi) = True	False : True =	6.4 : 1.0	False	<917>	912
has(honda) = True	False : True =	5.1 : 1.0	True	598	<1173>
has(toyota) = True	False : True =	4.9 : 1.0			
has(nexus) = True	False : True =	3.9 : 1.0			
has(manchester) = True	True : False =	3.7 : 1.0			
has(alcohol) = True	False : True =	3.7 : 1.0			
has(ufc160) = True	False : True =	3.7 : 1.0			
has(architecture) = True	False : True =	3.7 : 1.0			
has(barcelona) = True	True : False =	3.7 : 1.0			
has(iran) = True	True : False =	3.6 : 1.0			
has(snooki) = True	False : True =	3.3 : 1.0			

Features

- Features related to a particular tweet



- Features B: A + Tweet Properties
 - N. of Links
 - N. of Hash Tags

Results

F-measure: 0.61 - Accuracy: 0.62 - Precision: 0.62 - Recall: 0.60

			F	T
			a	r
			s	u
has(enviroment) = True	False : True =	16.3 : 1.0		
has(prometheus) = True	False : True =	11.1 : 1.0		
has(f1) = True	False : True =	9.1 : 1.0		
has(retweet) = True	True : False =	8.7 : 1.0		
has(sale) = True	False : True =	7.5 : 1.0		
has(mojito) = True	False : True =	7.3 : 1.0		
avg<3 = True	False : True =	7.1 : 1.0		
has(syria) = True	False : True =	6.8 : 1.0		
has(sopa) = True	False : True =	6.5 : 1.0		
has(mi) = True	False : True =	6.4 : 1.0	False	<1180> 650
has(honda) = True	False : True =	5.1 : 1.0	True	706<1064>
has(toyota) = True	False : True =	4.9 : 1.0		
has(nexus) = True	False : True =	3.9 : 1.0		
has(manchester) = True	True : False =	3.7 : 1.0		
has(architecture) = True	False : True =	3.7 : 1.0		
has(alcohol) = True	False : True =	3.7 : 1.0		
has(ufc160) = True	False : True =	3.7 : 1.0		
has(barcelona) = True	True : False =	3.7 : 1.0		
has(iran) = True	True : False =	3.6 : 1.0		
has(snooki) = True	False : True =	3.3 : 1.0		

Followers

The image shows two Twitter posts side-by-side, illustrating a comparison between them.

Post 1 (Top):

Sebastiano Milardo (@SebMilardo) - Happy
3:03 PM - 1 Giu 13

Risposta Elimina ★ Aggiungi ai preferiti Altro

Post 2 (Bottom):

Justin Bieber (@justinbieber) - Happy
6:28 PM - 31 Mag 13

Risposta Retweet ★ Aggiungi ai preferiti Altro

173.914 RETWEET 114.198 PREFERITI

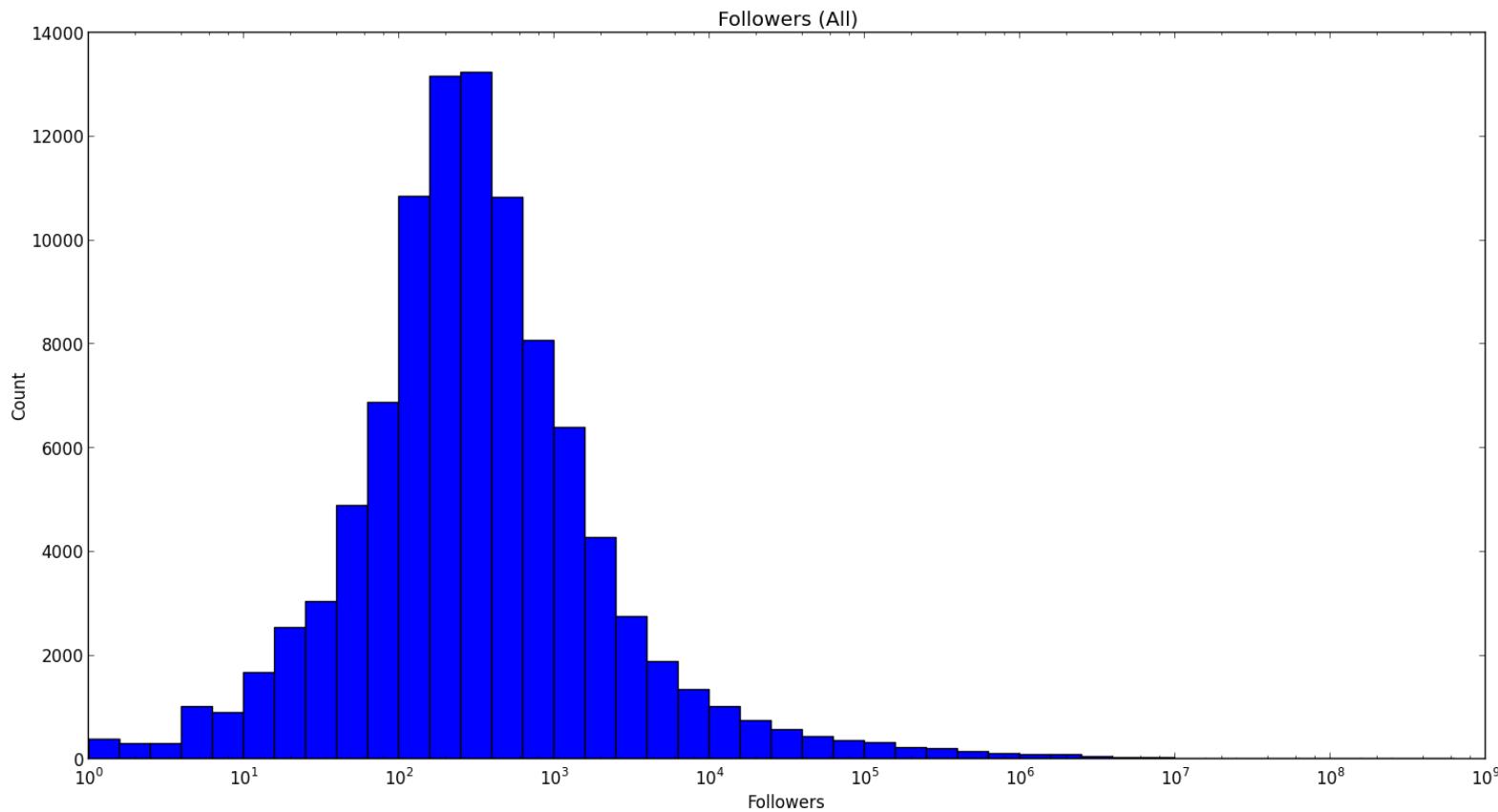
The posts are separated by a large, stylized "VS" graphic, indicating a competition or comparison between the two users' follower counts.

Top Retweeted

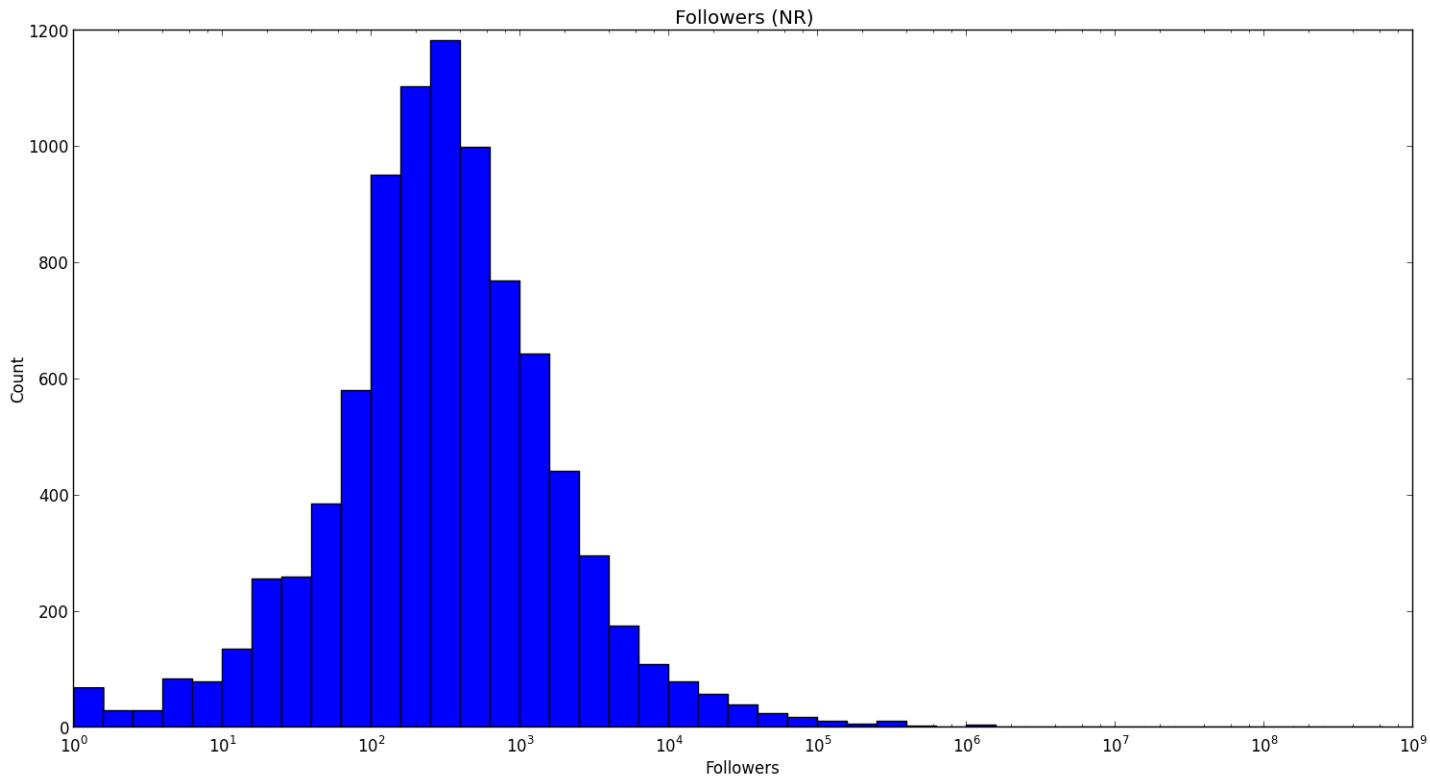
Count	Followers	Screen Name	Text
169266	10698555	Louis_Tomlinson	RT @Louis_Tomlinson: Always in my heart @Harry_Styles . Yours sincerely, Louis
157214	10627917	Real_Liam_Payne	RT @Real_Liam_Payne: Everybody meet mine and @daniellepeazer new dog Loki :) http://t.co/71ld7azgEW
126309	39610581	justinbieber	RT @justinbieber: My fans... My beliebers... are wild
118880	10698555	Louis_Tomlinson	RT @Louis_Tomlinson: What a great year together @EleanorJCalder :) Love you!!!!
112727	11620806	NiallOfficial	RT @NiallOfficial: Loki is chillin on my bunk! the dog is a legend! Saus monster http://t.co/Yyg9CmY6Nx
108165	10698555	Louis_Tomlinson	RT @Louis_Tomlinson: My heart misses DJ Malik @zaynmalik
71233	11620806	NiallOfficial	RT @NiallOfficial: Love always lads @jlsofficial http://t.co/e1GOSAvMls
66715	11620806	NiallOfficial	RT @NiallOfficial: Barcelona ! Muchas gracias !you were great! Soo loud! Xxx
61954	10698555	Louis_Tomlinson	RT @Louis_Tomlinson: Last night was unbelievable. I cannot believe we had the honour to perform at the Olympics! Has to be one of the best...
58188	39610581	justinbieber	RT @justinbieber: @pattiemalle
57242	13156420	Harry_Styles	RT @Harry_Styles: Cool last na
54535	11620806	NiallOfficial	RT @NiallOfficial: I can't bel
52254	11620806	NiallOfficial	RT @NiallOfficial: Morning Bar
51592	39610581	justinbieber	RT @justinbieber: seeing all t
49778	10698555	Louis_Tomlinson	RT @Louis_Tomlinson: Sad to he
49162	13156420	Harry_Styles	RT @Harry_Styles: Hapry Son. @
48931	13156420	Harry_Styles	RT @Harry_Styles: The first ti
48752	10627917	Real_Liam_Payne	RT @Real_Liam_Payne: Had a bri
45945	11620806	NiallOfficial	RT @NiallOfficial: I hate morn
44745	10627917	Real_Liam_Payne	RT @Real_Liam_Payne: Loving ba

Count	Followers	Screen Name	Text
169266	10698555	Louis_Tomlinson	RT @Louis_Tomlinson: Alway
157214	10627917	Real_Liam_Payne	RT @Real_Liam_Payne: Every
126309	39610581	justinbieber	RT @justinbieber: My fans.
118880	10698555	Louis_Tomlinson	RT @Louis_Tomlinson: What
112727	11620806	NiallOfficial	RT @NiallOfficial: Loki is
108165	10698555	Louis_Tomlinson	RT @Louis_Tomlinson: My he
71233	11620806	NiallOfficial	RT @NiallOfficial: Love al
66715	11620806	NiallOfficial	RT @NiallOfficial: Barcelo
61954	10698555	Louis_Tomlinson	RT @Louis_Tomlinson: Last
58188	39610581	justinbieber	RT @justinbieber: @pattiem
57242	13156420	Harry_Styles	RT @Harry_Styles: Cool las
54535	11620806	NiallOfficial	RT @NiallOfficial: I can't
52254	11620806	NiallOfficial	RT @NiallOfficial: Morning
51592	39610581	justinbieber	RT @justinbieber: seeing a
49778	10698555	Louis_Tomlinson	RT @Louis_Tomlinson: Sad t
49162	13156420	Harry_Styles	RT @Harry_Styles: Hapry So
48931	13156420	Harry_Styles	RT @Harry_Styles: The firs
48752	10627917	Real_Liam_Payne	RT @Real_Liam_Payne: Had a
45945	11620806	NiallOfficial	RT @NiallOfficial: I hate
44745	10627917	Real_Liam_Payne	RT @Real_Liam_Payne: Lovin

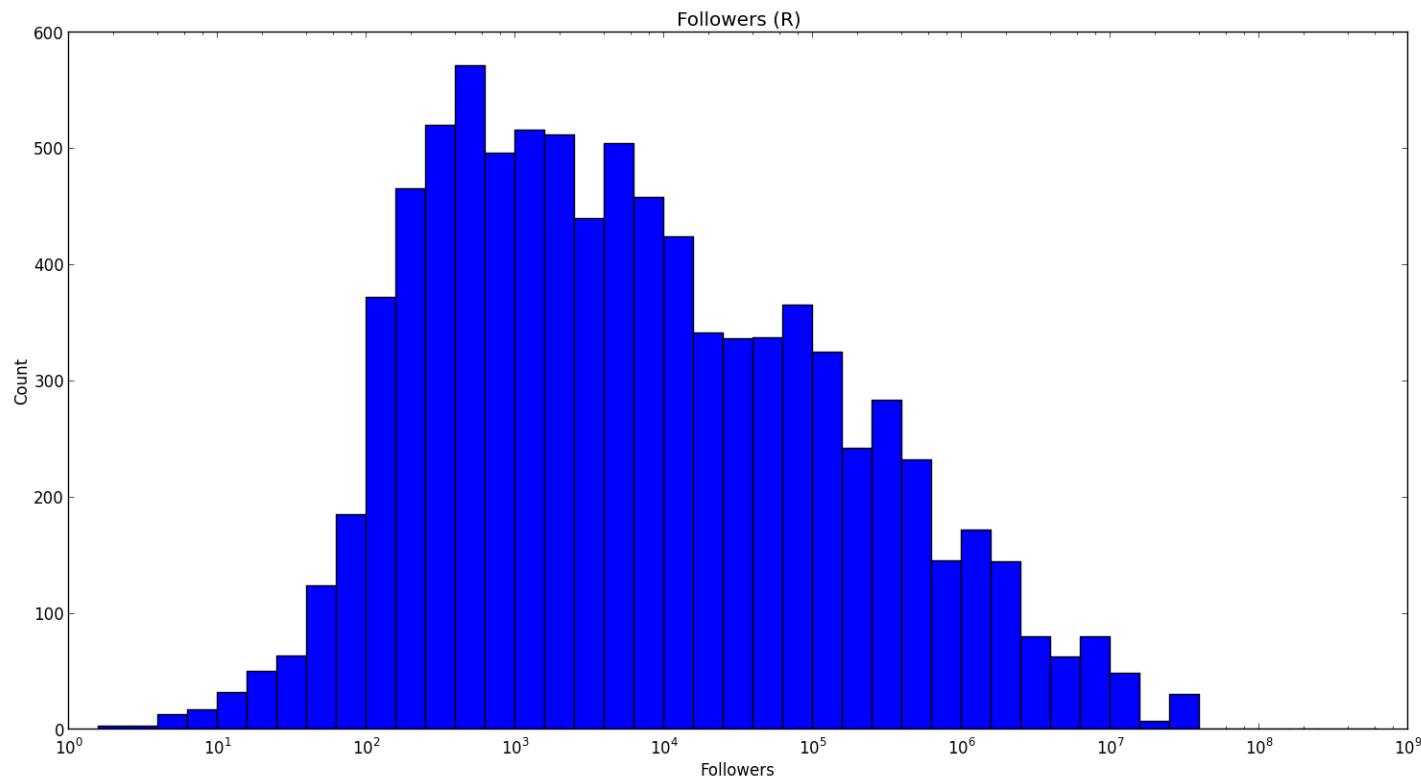
Followers (all tweets)



Followers (non retw.)



Followers (retw.)



Features

- Features related to a particular user.
- Features C: B + User Properties
 - N. Followers
 - N. Friends
 - N. Statuses
 - Account age
 - Verified
 - Default Profile



A Twitter profile card for Luca Parmitano. The card has a dark background. At the top right is a small square profile picture of him in a white spacesuit, giving a thumbs-up. Below the picture, his name "Luca Parmitano" is displayed in white, with a blue checkmark indicating verification. Underneath his name is his handle "@astro_luca". A bio follows: "I'm a European astronaut of Italian nationality. Italian Air Force test pilot. Onboard the ISS for Expedition 36/37. Volare!" At the bottom of the card, there are three metrics: "1.698 TWEET", "53 FOLLOWING", and "29.095 FOLLOWER". To the right of these metrics are two buttons: a grey "Following" button with a person icon, and a blue "Follow" button with a person icon and a plus sign.

Results

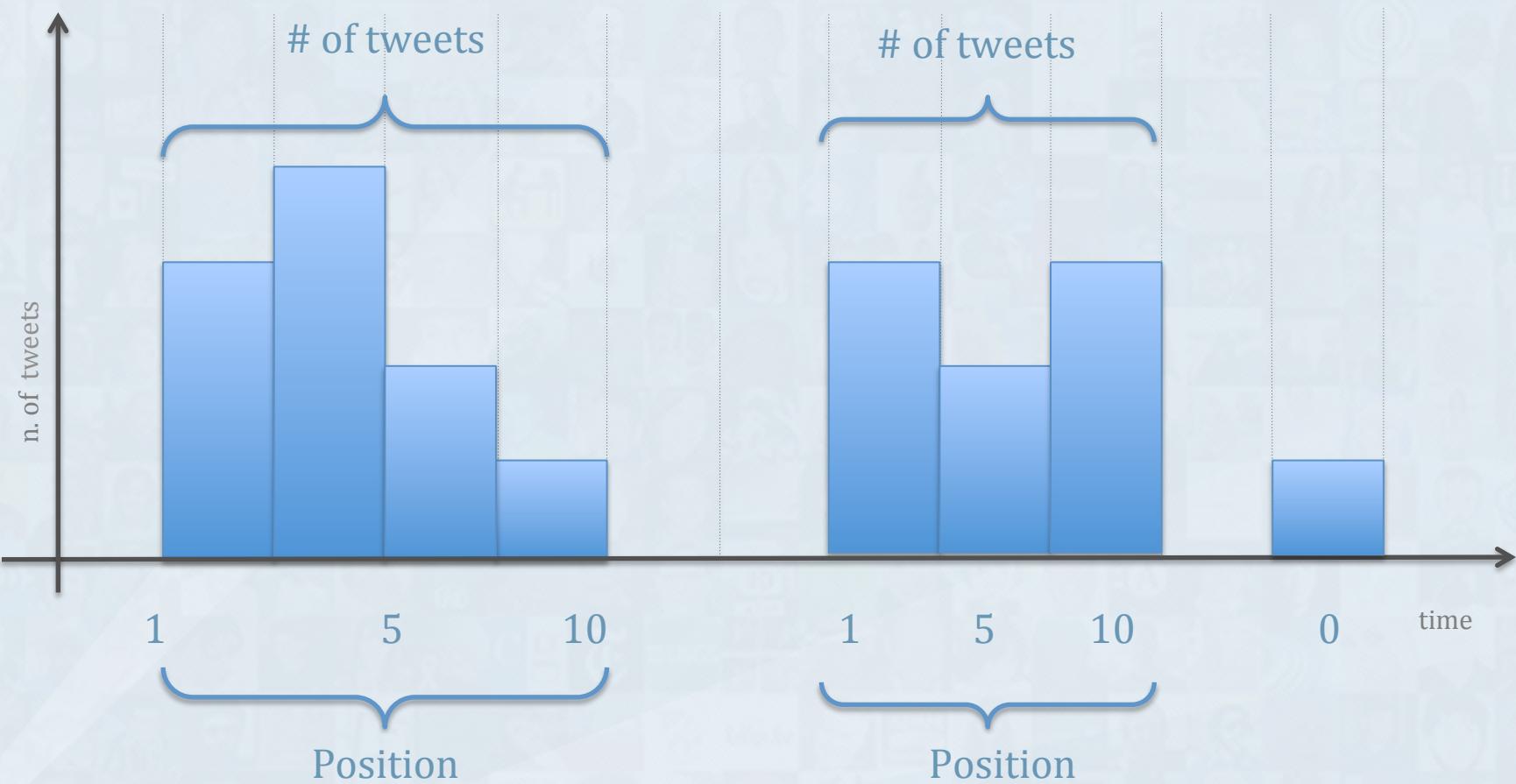
F-measure: 0.71 - Accuracy: 0.73 - Precision: 0.75 - Recall: 0.68

			F	T
			a	r
			s	u
5.0<=log_flw<7 = True	True : False =	40.1 : 1.0		
5.0<=log_frnds<7 = True	True : False =	18.3 : 1.0		
verified = True	True : False =	18.2 : 1.0		
has(enviroment) = True	False : True =	16.3 : 1.0		
log_flw<1.0 = True	False : True =	15.8 : 1.0		
has(prometheus) = True	False : True =	11.1 : 1.0		
has(f1) = True	False : True =	9.1 : 1.0		
has(retweet) = True	True : False =	8.7 : 1.0		
has(sale) = True	False : True =	7.5 : 1.0		
has(mojito) = True	False : True =	7.3 : 1.0	False	<1421> 408
avg<3 = True	False : True =	7.1 : 1.0	True	567<1204>
has(syria) = True	False : True =	6.8 : 1.0		
has(sopa) = True	False : True =	6.5 : 1.0		
has(mi) = True	False : True =	6.4 : 1.0		
has(honda) = True	False : True =	5.1 : 1.0		
has(toyota) = True	False : True =	4.9 : 1.0		
has(nexus) = True	False : True =	3.9 : 1.0		
has(manchester) = True	True : False =	3.7 : 1.0		
has(architecture) = True	False : True =	3.7 : 1.0		
has(alcohol) = True	False : True =	3.7 : 1.0		

Features

- There is something missing...

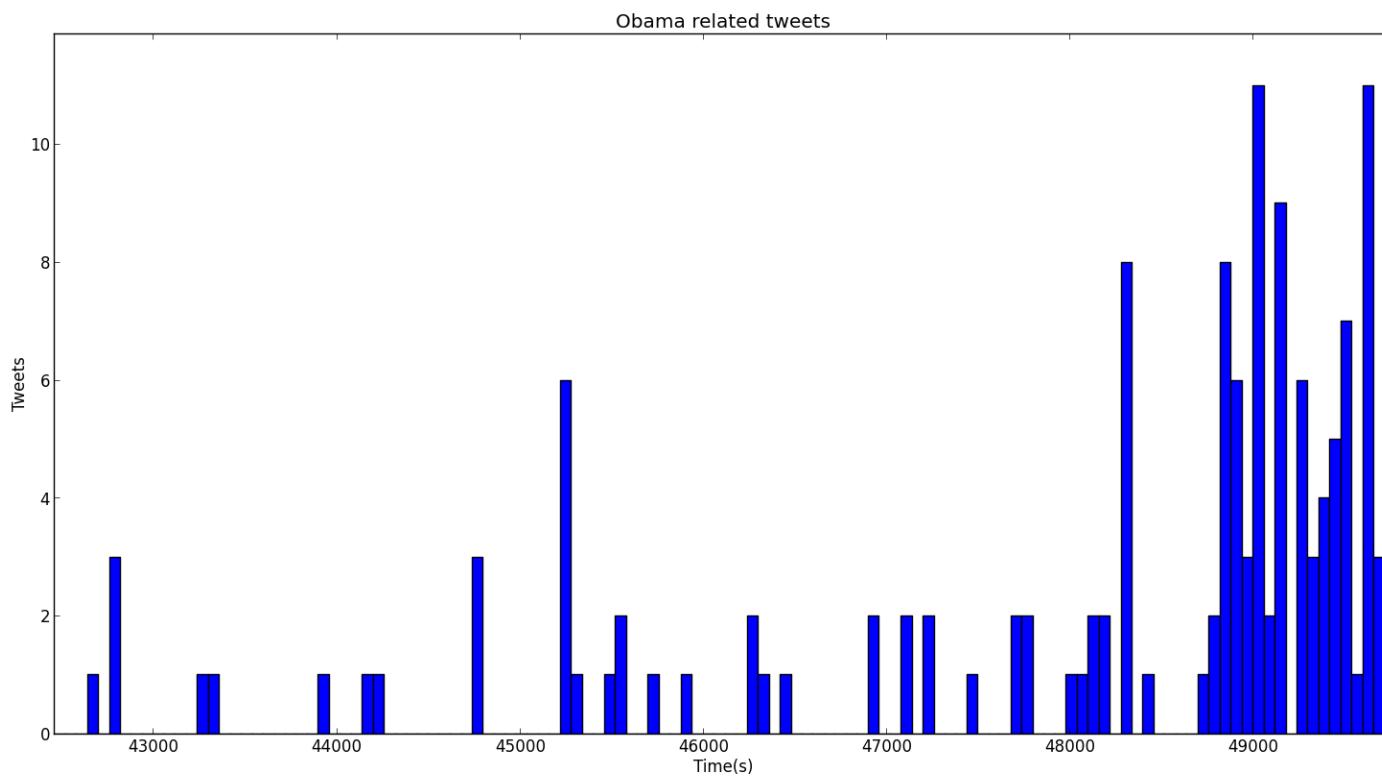
Time



Features

1369425527.0	0	0	Time (Unix Time)
1369479072.0	0	0	
1369485446.0	1.0	5	
1369485553.0	3.25	5	Start
1369485829.0	5.5	5	
1369486150.0	7.75	5	
1369486388.0	10.0	5	
1369512154.0	0	0	Stop
1369532095.0	0	0	
1369539965.0	1.0	3	
1369540138.0	5.5	3	
1369540254.0	10.0	3	
1369540866.0	0	0	
1369544351.0	0	0	
1369553350.0	0	0	

Time



Features

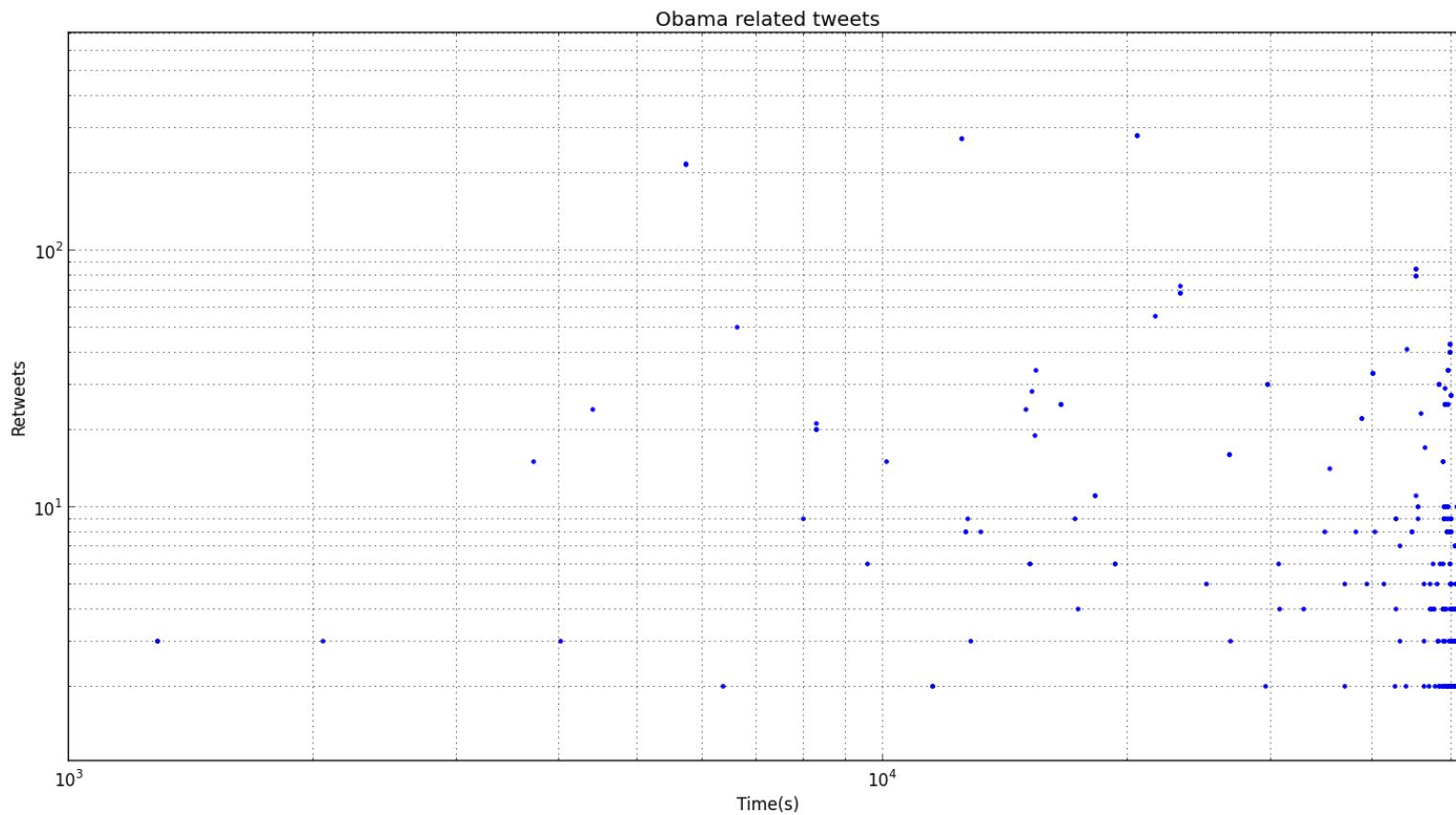
- Features related to time
- Features D: C + Time Properties
 - Is the first tweet in a series
 - Is the last tweet in a series
 - Is not in a series
 - Position in the series
 - Total number of Tweets in a series

Results

F-measure: 0.86 - Accuracy: 0.87 - Precision: 0.90 - Recall: 0.83

			F	T
			a	r
			s	u
s.pos==0 = True	True : False =	40.6 : 1.0		
5.0<=log_flw<7 = True	True : False =	40.1 : 1.0		
5.0<=log_frnds<7 = True	True : False =	18.3 : 1.0		
verified = True	True : False =	18.2 : 1.0		
has(enviroment) = True	False : True =	16.3 : 1.0		
log_flw<1.0 = True	False : True =	15.8 : 1.0		
has(prometheus) = True	False : True =	11.1 : 1.0		
s.pos==1 = True	True : False =	10.9 : 1.0		
has(f1) = True	False : True =	9.1 : 1.0		
has(retweet) = True	True : False =	8.7 : 1.0	False	<1657> 172
has(sale) = True	False : True =	7.5 : 1.0	True	294<1477>
has(mojito) = True	False : True =	7.3 : 1.0		
7.5< s.pos.< 10 = True	False : True =	7.2 : 1.0		
avg<3 = True	False : True =	7.1 : 1.0		
s.pos== 10 = True	True : False =	6.8 : 1.0		
has(syria) = True	False : True =	6.8 : 1.0		
has(sopa) = True	False : True =	6.5 : 1.0		
has(mi) = True	False : True =	6.4 : 1.0		
s.length<300.0 = True	True : False =	5.3 : 1.0		
has(honda) = True	False : True =	5.1 : 1.0		

...too good to be true



Results

F-measure: 0.85 - Accuracy: 0.86 - Precision: 0.87 - Recall: 0.83

			F	T
			a	r
			s	u
5.0<=log_flw<7 = True	True : False =	40.1 : 1.0		
5.0<=log_frnd<7 = True	True : False =	18.3 : 1.0		
verified = True	True : False =	18.2 : 1.0		
has(enviroment) = True	False : True =	16.3 : 1.0		
log_follow<1.0 = True	False : True =	15.8 : 1.0		
has(prometheus) = True	False : True =	11.1 : 1.0		
s.pos.==1 = True	True : False =	10.9 : 1.0		
has(f1) = True	False : True =	9.1 : 1.0		
has(retweet) = True	True : False =	8.7 : 1.0		
has(sale) = True	False : True =	7.5 : 1.0	False	<1607> 222
has(mojito) = True	False : True =	7.3 : 1.0	True	296<1475>
7.5<s.pos.< 10 = True	False : True =	7.2 : 1.0		
avg<3 = True	False : True =	7.1 : 1.0		
s.pos.== 10 = True	True : False =	6.8 : 1.0		
has(syria) = True	False : True =	6.8 : 1.0		
has(sopa) = True	False : True =	6.5 : 1.0		
has(mi) = True	False : True =	6.4 : 1.0		
s.length<300.0 = True	True : False =	5.3 : 1.0		
has(honda) = True	False : True =	5.1 : 1.0		
has(toyota) = True	False : True =	4.9 : 1.0		

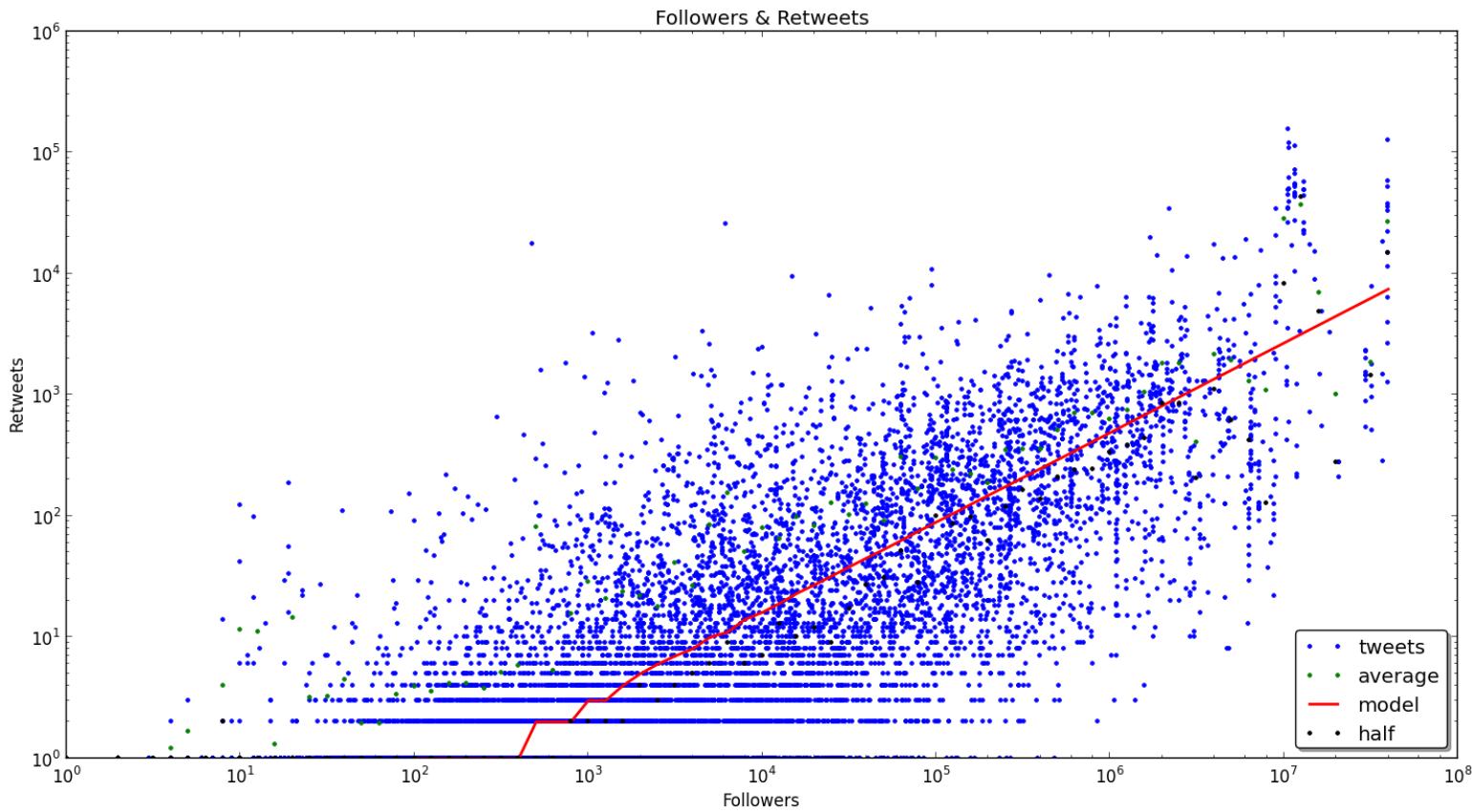
Comparison

Features Set	F-Measure	Accuracy	Precision	Recall
A	0,61	0,58	0,56	0,66
B	0,61	0,62	0,62	0,60
C	0,71	0,73	0,75	0,68
D	0,86	0,87	0,90	0,83
*D	0,85	0,86	0,87	0,83

Problem Description

- *What characterizes a tweet that gets retweeted by many people?*

Results



Retweet vs Followers

Expected retweet count = $a^{\log(f)-b}$

f = n. of followers

a in [5;6]

b in [2;2.5]

r_score = Retweet count/Expected retweet count

Examples:

$$5.5^{\log(f)-2.37}=1 \quad f=10$$

$$5.5^{\log(f)-2.37}=100 \quad f=160k$$

$$5.5^{\log(f)-2.37}=800\ 000 \quad f=31M$$

Retweet vs Followers

Was it a golden tweet? ($r_score = 1$? Maybe not...)

Barack Obama

@BarackObama

This account is run by Organizing for Action staff. Tweets from the President are signed -bo.

Washington, DC · barackobama.com

9.182 TWEET 662.115 FOLLOWING 32.238.024 FOLLOWER

Segui

Retweet vs Followers

- 4319 Good Retweets ($r_score > 1$)
- 4681 Bad Retweets ($r_score \leq 1$)

+

Bayesian Magic...

Comparison

Features Set	F-Measure	Accuracy	Precision	Recall
A	0,49	0,57	0,56	0,44
B	0,63	0,57	0,53	0,77
C	0,62	0,67	0,68	0,58
D	0,65	0,69	0,69	0,61
*D	0,64	0,68	0,69	0,59

Results

F-measure: 0.64 - Accuracy: 0.68 - Precision: 0.69 - Recall: 0.59

```
user<30Days = True  
has(goodmorning) = True  
    has(missing) = True  
5.0<=log_statuses<7 = True  
1.0<=log_statuses<3.0 = True  
    has(senate) = True  
    has(retweet) = True  
    has(oil) = True  
1Month<=user<6Months = True  
    has(monaco) = True  
    has(snooki) = True  
    has(bcci) = True  
    has(fail) = True  
    has(uk) = True  
    has(ufc) = True  
    has(war) = True  
3.0<=log_statuses<5.0 = False  
    has(lumia) = True  
    has(union) = True  
    has(boxing) = True
```

True : False =	12.2 : 1.0	
True : False =	7.1 : 1.0	
False : True =	6.0 : 1.0	
False : True =	5.9 : 1.0	
True : False =	4.8 : 1.0	
False : True =	4.7 : 1.0	
True : False =	4.2 : 1.0	-----
False : True =	4.0 : 1.0	False <726>226
True : False =	3.7 : 1.0	True 348<500>
False : True =	3.5 : 1.0	-----
False : True =	3.5 : 1.0	
False : True =	3.4 : 1.0	
True : False =	3.2 : 1.0	
False : True =	3.1 : 1.0	
False : True =	3.1 : 1.0	
False : True =	3.1 : 1.0	
True : False =	3.0 : 1.0	
False : True =	2.8 : 1.0	
False : True =	2.8 : 1.0	
False : True =	2.8 : 1.0	

F	a	T
l	r	
s	u	
e	e	
-----	+-----+	
False	<726>226	
True	348<500>	
-----	+-----+	

Conclusions

- Recipe for a golden tweet:
 - Popularity is the key!
 - Try to be the first one... or the last one
 - At the beginning everything is easier
 - Don't write too many statuses
 - Ask for a retweet

Questions



Thank You!