

## DB Report

The table provided was not normalized and having following columns

created\_at, text, tweet\_id, in\_reply\_to\_screen\_name, in\_reply\_to\_status\_id, in\_reply\_to\_user\_id, retweet\_count, tweet\_source, retweet\_of\_tweet\_id, hashtag1, hashtag2, hashtag3, hashtag4, hashtag5, hashtag6, user\_id, user\_name, user\_screen\_name, user\_location, user\_utc\_offset, user\_time\_zone, user\_followers\_count, user\_friends\_count, user\_lang, user\_description, user\_status\_count, user\_created\_at

There are total five subject areas:

1. Tweet\_User
2. Tweet\_Reply
3. Tweet\_Retweet
4. Tweet\_source
5. Tweet\_hashtag

For each subject area is required individual table as per normalization rule. To remove redundancy.

However, we can create a table for language but in this situation is not required.

After performing 3<sup>rd</sup> Normalized from tables which will populate looks like this:

## **Table Tweet\_User**

```
(  
  user_id,  
  user_name  
  user_screen_name  
  user_location  
  user_utc_offset,  
  user_time_zone  
  user_followers_count,  
  user_friends_count,  
  user_lang  
  user_description  
  user_status_count,  
  user_created_at  
)
```

## **TABLE Tweet\_Reply**

```
(  
  tweet_Reply_id,  
  tweet_id ,  
  in_reply_to_screen_name,  
  in_reply_to_status_id,  
  in_reply_to_user_id  
)
```

### **Table Tweet\_Retweet**

```
( retweet_id,  
  retweet_of_tweet_id ,  
  retweet_count,  
  tweet_id  
)
```

### **TABLE Tweet\_source**

```
(  
  tweet_id ,  
  tweet_source,  
  text,  
  retweet_id,  
  tweet_Reply_id,  
  user_id    ,  
  created_at  
)
```

### **TABLE Tweet\_hashtag**

```
(  
  hashtagid,  
  hashtag,  
  tweet_id
```

)

### **Functional Dependencies:**

user\_id -> user\_name , user\_screen\_name

user\_lang , user\_created\_at -> user\_id

user\_time\_zone -> user\_id

tweet\_Reply\_id -> tweet\_id

retweet\_id -> retweet\_of\_tweet\_id

hashtagid-> hashtag

Note: Before loading data into final tables I've removed duplicated while populating data into new tables.

Copy the data provided:

Query EditorQuery History

Scratch Pad

```
1 COPY bad_giant_table
2 -- (created_at, text, tweet_id, in_reply_to_screen_name, in_reply_to_status_id, in_reply_to_user_id, retweet_count, tweet_source)
3 FROM 'F:\bad_giant_table.csv'
4 DELIMITER ','
5 CSV HEADER;
6
7
```

Data Output

Explain

Messages

Notifications

COPY 110573

Query returned successfully in 17 secs 141 msec.

```
1 SELECT created_at, text, tweet_id, in_reply_to_screen_name, in_reply_to_status_id, in_reply_to_user_id, retweet_count, tweet_source
2 FROM public.bad_shorter_table;
```

Data Output

Explain

Messages

Notifications

	created_at timestamp with time zone	text character varying (255)	tweet_id [PK] bigint	in_reply_to_screen_name character varying (255)	in_reply_to_status_id bigint	in_reply_to_user_id bigint	retweet_count integer	tweet_source character varying (255)	retweet_of_tweet_id bigint
1	2012-10-22 08:00:00+05	Everybody cant be spitting rea...	13890351239168	[null]		[null]		0 <a href="http://twitter.com/do...	[null]
2	2012-10-22 08:00:00+05	Watching Bad Boys 2 love this...	13890363817984	[null]		[null]		0 <a href="http://www.ig.com rel...	[null]
3	2012-10-22 08:00:00+05	Balik KL for a day before head...	13890355445762	[null]		[null]		0 <a href="http://seesmic.com/ r...	[null]
4	2012-10-22 08:00:00+05	Too much on my mind cant sl...	13890355429376	[null]		[null]		0 <a href="http://twitter.com/do...	[null]
5	2012-10-22 08:00:00+05	Just thinking thoe , cause i ml...	13890372210688	[null]		[null]		0 <a href="http://twitter.com/do...	[null]
6	2012-10-22 08:00:00+05	tava conversando com a julia ...	13890368012289	[null]		[null]		0 <a href="http://twitter.com/do...	[null]
7	2012-10-22 08:00:00+05	"@LHopkins07: #NationalFear...	13890372218880	[null]		[null]		0 <a href="http://twitter.com/do...	[null]
8	2012-10-22 08:00:00+05	"@FUN: If a mosquito is biting...	13890372231170	[null]		[null]		0 <a href="http://twitter.com/do...	[null]
9	2012-10-22 08:00:00+05	Mereka mirip ... RT @EXOSun...	13890380619776	[null]		[null]		0 <a href="http://bit.ly/fVzRUd re...	[null]
10	2012-10-22 08:00:00+05	RT @ohMG91: I'll save tomorr...	13890359640064	[null]		[null]		0 <a href="http://twitter.com/do...	[null]
11	2012-10-22 08:00:00+05	@Ivania_Gabby im sorry boob...	13890384809984	Ivania_Gabby	260213647538806784	411004777	0	web	[null]
12	2012-10-22 08:00:00+05	And... its gone!	13890363826176	[null]		[null]		0 <a href="http://www.facebook...	[null]
13	2012-10-22 08:00:00+05	@Shelby_St_James: Gonna se...	13890355453954	[null]		[null]		0 <a href="http://twitter.com/dev...	[null]
14	2012-10-22 08:00:00+05	@haziqzeke eh meet me laaa...	13890368040961	haziqzeke	260213610607955968	44981750	0	<a href="http://www.echofon.c...	[null]

```
1 SELECT created_at, text, tweet_id, in_reply_to_screen_name, in_reply_to_status_id, in_reply_to_user_id, retweet_count, tweet_source
2 FROM public.bad_giant_table;
```

Data Output

Explain

Messages

Notifications

	created_at timestamp with time zone	text character varying (255)	tweet_id [PK] bigint	in_reply_to_screen_name character varying (255)	in_reply_to_status_id bigint	in_reply_to_user_id bigint	retweet_count integer	tweet_source character varying (255)	retweet_of_tweet_id bigint
1	2012-10-22 08:00:57+05	0101 @Idorminhoca te amo !	14129447563264	[null]		[null]		0 web	
2	2012-10-22 08:00:58+05	@Feftherrero nunca -	14133658619905	Feftherrero	260212715052756993	165989321	0	web	
3	2012-10-22 08:01:09+05	Voyyy RT @FeelLikeAKaren: S...	14179766620162	[null]		[null]		0 web	
4	2012-10-22 08:01:14+05	#GeorgeStrait	14200767500288	[null]		[null]		0 web	
5	2012-10-22 08:00:00+05	Start looking towards the sky .	13889281708032	[null]		[null]		0 web	
6	2012-10-22 08:01:29+05	Cool.	14263665274881	[null]		[null]		0 web	
7	2012-10-22 08:00:50+05	Its not getting easier	14099873505280	[null]		[null]		0 web	
8	2012-10-22 08:01:50+05	Un suspiro más y lo aspiro.	14351745671168	[null]		[null]		0 web	
9	2012-10-22 08:01:49+05	HANNAH HAS 800 FOLLOWE...	14347563933697	[null]		[null]		0 web	
10	2012-10-22 08:01:53+05	hush, hush.	14364307603456	[null]		[null]		0 web	
11	2012-10-22 08:01:53+05	andre 3000 kill errr thing	14364328558593	[null]		[null]		0 web	
12	2012-10-22 08:01:57+05	editando. o.o	14381114171393	[null]		[null]		0 web	
13	2012-10-22 08:02:03+05	RT @Cydiapomelia: http://t.co...	14406271614976	[null]		[null]		4 <a href="http://blackberry.co...	26021271634459
14	2012-10-22 08:02:08+05	0101 meu amor ♥	14427259920384	[null]		[null]		0 web	
15	2012-10-22 08:02:09+05	ya wassalam	14431437422592	[null]		[null]		0 <a href="http://twitter.com/de...	

## SQL Queries:

1

```
4 a) How many tweets are there in total?
5
6 Select
7 count(*) as total_tweet
8 from Tweet_source;
9
```

Data Output

Explain

Messages

Notifications

	<div>total_tweet</div> <div>bigint</div>	
1	43876	

2

10 b) How are these tweets distributed across languages? Write a query that shows,  
11 for every language ( user\_lang ) the number of tweets in that language.

```
12  
13 Select  
14 user_lang,  
15 count(*) as total_tweet_Per_Language  
16 from Tweet_user  
17 group by 1;
```

Data Output Explain Messages Notifications

	user_lang character varying (10)	total_tweet_per_language bigint
1	fr	221
2	tr	121
3	en	69279
4	fi	1
5	ru	382
6	cs	2
7	ja	9596
8	sv	2
9	fil	9
10	eu	1
11	de	71
12	nl	60
13	id	1360
14	zh-cn	21
15	zh-tw	14
16	ar	1328
17	no	1
18	th	217





4

39 a) What fraction of the tweets are retweets ?

40

41 **Select**

42 **t tweet\_id,**

43 **(count( rt.tweet\_id) \* 1.0/ (select count(\*) from Tweet\_source) ) as Fraction\_Tweet**

44 **from Tweet\_source t**

45 **inner join Tweet\_Retweet rt**

46 **on t.retweet\_id = rt.retweet\_id**

47 **group by 1**

48

49

50

Data Output Explain Messages Notifications

	tweet_id [PK] bigint	fraction_tweet numeric
1	260214557283336192	0.000022791503327559485824
2	260219334582992896	0.000022791503327559485824
3	260216721514852352	0.000022791503327559485824
4	260219468796530688	0.000022791503327559485824
5	260221456879857665	0.000022791503327559485824
6	260217409393274880	0.000022791503327559485824
7	260217124168036352	0.000022791503327559485824
8	260220852929433600	0.000022791503327559485824
9	260215245136621568	0.000022791503327559485824
10	260219997274656768	0.000022791503327559485824
11	260217317135360000	0.000022791503327559485824
12	260214066545586176	0.000022791503327559485824
13	260217702973571072	0.000022791503327559485824
14	260214725038727168	0.000022791503327559485824
15	260219540091326464	0.000022791503327559485824
16	260214804734685188	0.000022791503327559485824

51 b) Compute the average **number of** retweets per tweet.

52

53

54 **Select**

55 **t.tweet\_id,**

56 **avg( rt.retweet\_count) as Avg\_Retweet**

57 **from Tweet\_source t**

58 **inner join Tweet\_Retweet rt**

59 **on t.retweet\_id = rt.retweet\_id**

60 **group by 1**

61

62

Data Output Explain Messages Notifications

	tweet_id [PK] bigint	avg_retweet numeric	
1	260214557283336192	1.00000000000000000000	
2	260216721514852352	1.00000000000000000000	
3	260219468796530688	1.00000000000000000000	
4	260221456879857665	0.00000000000000000000	
5	260217409393274880	1.00000000000000000000	
6	260217124168036352	2.00000000000000000000	
7	260220852929433600	9.00000000000000000000	
8	260215245136621568	1.00000000000000000000	
9	260219997274656768	1.00000000000000000000	
10	260217317135360000	1.00000000000000000000	
11	260214066545586176	0.00000000000000000000	
12	260217702973571072	1.00000000000000000000	
13	260214725038727168	0.00000000000000000000	
14	260219540091326464	2.00000000000000000000	
15	260214804734685188	10.00000000000000000000	
16	260214544700432384	3.00000000000000000000	

6

64 c) What fraction of the tweets are never retweeted?

65

66 **Select**

67 **t.tweet\_id,**

68 **(count( t.tweet\_id) \* 1.0/ (select count(\*) from Tweet\_source) ) as Fraction\_Tweet**

69 **from Tweet\_source t**

70 **left join Tweet\_Retweet rt**

71 **on t.retweet\_id = rt.retweet\_id**



72 **where t.retweet\_id is null**

73 **group by 1**

74

75

Data Output Explain Messages Notifications

	 tweet_id [PK] bigint	 fraction_tweet numeric	
1	260215006052876288	0.000022791503327559485824	
2	260217652667092993	0.000022791503327559485824	
3	260216780260274176	0.000022791503327559485824	
4	260215683760140291	0.000022791503327559485824	
5	260214263677845504	0.000022791503327559485824	
6	260214435635933184	0.000022791503327559485824	
7	260219183583866880	0.000022791503327559485824	
8	260220056015884288	0.000022791503327559485824	
9	260219734212091904	0.000022791503327559485824	
10	260214653714587648	0.000022791503327559485824	
11	260219686917128193	0.000022791503327559485824	
12	260217921098362880	0.000022791503327559485824	
13	260215064773160960	0.000022791503327559485824	
14	260215127704494080	0.000022791503327559485824	
15	260221301724168192	0.000022791503327559485824	

76 d) What fraction of the tweets are retweeted fewer times than the average number of retweets (and what does this say about the distribution)?

```

77
78 select TT.tweet_id,
79 TT.Fraction_Tweet,
80 TT.Avg_Retweet from
81 (
82   Select
83     t.tweet_id,
84     (count( rt.tweet_id) * 1.0/ (select count(*) from Tweet_source) ) as Fraction_Tweet,
85     avg( rt.retweet_count) as Avg_Retweet
86   from Tweet_source t
87   inner join Tweet_Retweet rt
88   on t.retweet_id = rt.retweet_id
89   group by 1
90 ) TT
91 where TT.Fraction_Tweet < TT.Avg_Retweet
92
93

```

Data Output Explain Messages Notifications

	tweet_id [PK] bigint	fraction_tweet numeric	avg_retweet numeric	
1	260213890351239169	0.000022791503327559485824	106.000000000000000000	
2	260213890351239170	0.000022791503327559485824	1.00000000000000000000	
3	260213890351255552	0.000022791503327559485824	5027.0000000000000000	
4	260213890355453952	0.000022791503327559485824	112.0000000000000000	
5	260213890359635968	0.000022791503327559485824	1.000000000000000000	
6	260213890368012291	0.000022791503327559485824	2.000000000000000000	
7	260213890368016384	0.000022791503327559485824	7.000000000000000000	
8	260213890368024576	0.000022791503327559485824	1.000000000000000000	
9	260213890368040962	0.000022791503327559485824	80.0000000000000000	
10	260213890372222976	0.000022791503327559485824	30.0000000000000000	

8

102 a) What is the number of distinct hashtags found in these tweets?

103

104

105 **Select distinct** hashtag

106 **from** Tweet\_hashtag

107

108

Data Output Explain Messages Notifications

	hashtag character varying (150)	
1	dontwakemeup	
2	spyair	
3	غادة_شبير	
4	IDX	
5	burkesquad	
6	stress	
7	believetour	
8	Tmit	
9	DOIT	
10	10FactsAboutMe	
11	cmnews	
12	GospelRap	
13	Madvanna	
14	Swacked	
15	nottacos	
16	Tyga	

110 b) What **are** the top ten most popular hashtags, **by number of** usages?

111

112

113 **Select** hashtag,

114 **count**( hashtag ) **as** Top\_Hashtags

115 **from** Tweet\_hashtag

116 **group by** 1

117 **order by** 2 desc limit 10

118

119

120 c) Write a query which, for each language, the top three most popular hashtags are

Data Output Explain Messages Notifications

	hashtag character varying (150)	top_hashtags bigint
1	ReasonsIFailAtBeingAGirl	467
2	RED	205
3	oomf	188
4	HonestyHour	172
5	TeamFollowBack	139
6	EresGuapaSi	130
7	10PeopleYouTrulyLove	126
8	TweetLikeAGirl	98
9	ImSingleBecause	97
10	WeAllGotThatOneFriend	96

10

```
120 c) Write a query giving, for each language, the top three most popular hashtags in that language.
121
122 Select hashtag,user_lang, count( hashtag )
123 from Tweet_hashtag
124 join tweet_source
125 on Tweet_hashtag.tweet_id = tweet_source.tweet_id
126 join tweet_user
127 on tweet_source.user_id = tweet_user.user_id
128 group by 1,2
129 order by 3 desc limit 3
130
```

Data Output Explain Messages Notifications

	hashtag character varying (150)	user_lang character varying (10)	count bigint
1	ReasonsIFailAtBeingAGirl	en	184
2	RED	en	84
3	MagicMike	en	78

11

```
139 a) How many tweets are neither replies, nor replied to?
140
141 select count(*) as count_Noreply
142 from tweet_source
143 where tweet_Reply_id is null
144
145
```

Data Output Explain Messages Notifications

	count_noreply bigint
1	21438

147 b) If a user user1 replies to another user user2 , what is the probability that they have the same language setting?

```
148
149 Select user_lang,
150 COUNT(user_lang)*1.0/(SELECT count(tweet_id) FROM tweet_source) as Prob_UselangSetting
151 from tweet_source as ts
152 JOIN tweet_reply as tu
153 on ts.tweet_Reply_id= tu.tweet_Reply_id
154 join tweet_user as tu1
155 on ts.user_id=tu1.user_id
156 group by 1
157
158
```

Data Output Explain Messages Notifications



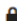
	user_lang character varying (10)	prob_uselangsetting numeric
1	ar	0.00423921961892606436
2	ca	0.000022791503327559485824
3	de	0.00027349803993071383
4	en	0.32728598778375421643
5	es	0.08601513355820949950
6	fil	0.000045583006655118971647
7	fr	0.00109399215972285532
8	id	0.00556112681192451454
9	it	0.00013674901996535691
10	ja	0.06212963807092715836
11	ko	0.00394293007566779105
12	msa	0.000045583006655118971647
13	nl	0.00025070653660315434
14	no	0.000022791503327559485824
15		



159 c) How does this compare to the probability that two arbitrary users have the same language setting?  
160 Throughout, you may create views that support your queries.

```
161  
162 WITH lang_setting AS(  
163   Select   ts.user_id,  
164   COUNT(user_lang)*1.0/(SELECT count(tweet_id) FROM tweet_source) as Prob_UseLangSetting  
165   from tweet_source as ts  
166   JOIN tweet_reply as tu  
167   on ts.tweet_reply_id= tu.tweet_reply_id  
168   join tweet_user as tu1  
169   on ts.user_id=tu1.user_id  
170   AND tu1.user_lang = tu1.user_lang  
171   group by 1  
172 )  
173 Select * from lang_setting ORDER BY RANDOM() LIMIT 2  
174
```

Data Output Explain Messages Notifications

	 user_id integer	 prob_uselangsetting numeric	
1	548028773	0.000022791503327559485824	
2	746909557	0.000022791503327559485824	