```
alter table May2015
drop column created utc, link id, score hidden, author flair css class,
author flair text,
removal reason, archived, retrieved on, edited, controversiality,
parent id;
--remove non default subreddits excluding announcements
delete from May2015
where not subreddit='Art'
and not subreddit='AskReddit'
and not subreddit='askscience'
and not subreddit='aww'
and not subreddit='blog'
and not subreddit='books'
and not subreddit='creepy'
and not subreddit='dataisbeautiful'
and not subreddit='DIY'
and not subreddit='Documentaries'
and not subreddit='EarthPorn'
and not subreddit='explainlikeimfive'
and not subreddit='food'
and not subreddit='funny'
and not subreddit='Futurology'
and not subreddit='gadgets'
and not subreddit='gaming'
and not subreddit='GetMotivated'
and not subreddit='gifs'
and not subreddit='history'
and not subreddit='IAmA'
and not subreddit='InternetIsBeautiful'
and not subreddit='Jokes'
and not subreddit='LifeProTips'
and not subreddit='listentothis'
and not subreddit='mildlyinteresting'
```

```
and not subreddit='movies'
and not subreddit='Music'
and not subreddit='news'
and not subreddit='nosleep'
and not subreddit='nottheonion'
and not subreddit='OldSchoolCool'
and not subreddit='personalfinance'
and not subreddit='philosophy'
and not subreddit='photoshopbattles'
and not subreddit='pics'
and not subreddit='science'
and not subreddit='Showerthoughts'
and not subreddit='space'
and not subreddit='sports'
and not subreddit='television'
and not subreddit='tifu'
and not subreddit='todayilearned'
and not subreddit='UpliftingNews'
and not subreddit='videos'
and not subreddit='worldnews';
select count(distinct subreddit) from May2015;
select distinct subreddit from May2015;
delete from May2015 where body='[deleted]';
scp may 2015 trim.csv dkang10@129.150.64.74:/home/dkang10
```

```
hdfs dfs -mkdir group3
hdfs dfs -mkdir group3/dataset
hdfs dfs -ls group3
--put file in hdfs
hdfs dfs -put may 2015 trim.csv group3/dataset
hdfs dfs -ls group3/dataset
--change permissions of new directory
hdfs dfs -chmod -R o+w .
create external table may 2015 (
   id string,
   subreddit id string,
   subreddit string,
   name string,
   author string,
   ups int,
   downs int,
   score int,
   gilted int,
   distinguished string,
   body string
row format delimited
fields terminated by '|'
lines terminated by '\n'
location '/user/dkang10/group3/dataset'
tblproperties('skip.header.line.count'='1');
--check new table
show tables;
describe may 2015;
select count(*) from may_2015;
```

```
--clean table from csv and text box errors
create table may 2015 a as
select id, subreddit id, subreddit, name, author, ups, downs, score,
gilted, distinguished, body from may 2015
where subreddit='Art'
or subreddit='AskReddit'
or subreddit='askscience'
or subreddit= 'aww'
or subreddit='blog'
or subreddit='books'
or subreddit='creepy'
or subreddit='dataisbeautiful'
or subreddit='DIY'
or subreddit='Documentaries'
or subreddit='EarthPorn'
or subreddit='explainlikeimfive'
or subreddit='food'
or subreddit='funny'
or subreddit='Futurology'
or subreddit='gadgets'
or subreddit='gaming'
or subreddit='GetMotivated'
or subreddit='gifs'
or subreddit='history'
or subreddit='IAmA'
or subreddit='InternetIsBeautiful'
or subreddit='Jokes'
or subreddit='LifeProTips'
or subreddit='listentothis'
or subreddit='mildlyinteresting'
or subreddit='movies'
or subreddit='Music'
or subreddit='news'
or subreddit='nosleep'
or subreddit='nottheonion'
or subreddit='OldSchoolCool'
or subreddit='personalfinance'
or subreddit='philosophy'
```

```
or subreddit='photoshopbattles'
or subreddit='pics'
or subreddit='science'
or subreddit='Showerthoughts'
or subreddit='space'
or subreddit='sports'
or subreddit='television'
or subreddit='tifu'
or subreddit='todayilearned'
or subreddit='UpliftingNews'
or subreddit='videos'
or subreddit='worldnews';
--general analysis begins~~~~~~~~~~~~~~~~~~~
--query comments with highest upvotes across default subreddits
select subreddit, ups, downs, score, body from may 2015 a
order by score desc limit 25;
   statement | 1902 | 1 | 1902 | 1 | 1903 | 1 | 1904 | 1 | 1905 | 1 | 1905 | 1 | 1905 | 1 | 1905 | 1 | 1905 | 1 | 1905 | 1 | 1905 | 1 | 1905 | 1 | 1905 | 1 | 1905 | 1 | 1905 | 1 | 1905 | 1 | 1905 | 1 | 1905 | 1 | 1905 | 1 | 1905 | 1 | 1905 | 1 | 1905 | 1 | 1905 | 1 | 1905 | 1 | 1905 | 1 | 1905 | 1 | 1905 | 1 | 1905 | 1 | 1905 | 1 | 1905 | 1 | 1905 | 1 | 1905 | 1 | 1905 | 1 | 1905 | 1 | 1905 | 1 | 1905 | 1 | 1905 | 1 | 1905 | 1 | 1905 | 1 | 1905 | 1 | 1905 | 1 | 1905 | 1 | 1905 | 1 | 1905 | 1 | 1905 | 1 | 1905 | 1 | 1905 | 1 | 1905 | 1 | 1905 | 1 | 1905 | 1 | 1905 | 1 | 1905 | 1 | 1905 | 1 | 1905 | 1 | 1905 | 1 | 1905 | 1 | 1905 | 1 | 1905 | 1 | 1905 | 1 | 1905 | 1 | 1905 | 1 | 1905 | 1 | 1905 | 1 | 1905 | 1 | 1905 | 1 | 1905 | 1 | 1905 | 1 | 1905 | 1 | 1905 | 1 | 1905 | 1 | 1905 | 1 | 1905 | 1 | 1905 | 1 | 1905 | 1 | 1905 | 1 | 1905 | 1 | 1905 | 1 | 1905 | 1 | 1905 | 1 | 1905 | 1 | 1905 | 1 | 1905 | 1 | 1905 | 1 | 1905 | 1 | 1905 | 1 | 1905 | 1 | 1905 | 1 | 1905 | 1 | 1905 | 1 | 1905 | 1 | 1905 | 1 | 1905 | 1 | 1905 | 1 | 1905 | 1 | 1905 | 1 | 1905 | 1 | 1905 | 1 | 1905 | 1 | 1905 | 1 | 1905 | 1 | 1905 | 1 | 1905 | 1 | 1905 | 1 | 1905 | 1 | 1905 | 1 | 1905 | 1 | 1905 | 1 | 1905 | 1 | 1905 | 1 | 1905 | 1 | 1905 | 1 | 1905 | 1 | 1905 | 1 | 1905 | 1 | 1905 | 1 | 1905 | 1 | 1905 | 1 | 1905 | 1 | 1905 | 1 | 1905 | 1 | 1905 | 1 | 1905 | 1 | 1905 | 1 | 1905 | 1 | 1905 | 1 | 1905 | 1 | 1905 | 1 | 1905 | 1 | 1905 | 1 | 1905 | 1 | 1905 | 1 | 1905 | 1 | 1905 | 1 | 1905 | 1 | 1905 | 1 | 1905 | 1 | 1905 | 1 | 1905 | 1 | 1905 | 1 | 1905 | 1 | 1905 | 1 | 1905 | 1 | 1905 | 1 | 1905 | 1 | 1905 | 1 | 1905 | 1 | 1905 | 1 | 1905 | 1 | 1905 | 1 | 1905 | 1 | 1905 | 1 | 1905 | 1 | 1905 | 1 | 1905 | 1 | 1905 | 1 | 1905 | 1 | 1905 | 1 | 1905 | 1 | 1905 | 1 | 1905 | 1 | 1905 | 1 | 1905 | 1 | 1905 | 1 | 1905 | 1 | 1905 | 1 | 1905 | 1 | 1905 | 1 | 1905 | 1 | 1905 | 1 | 1905 | 1 | 1905 | 1 | 1905 | 1 | 1905 | 1 | 1905 | 1 | 1905 | 1 | 1905 | 1 | 1905 | 1 | 1905 | 1 | 1905 | 1 | 1905 | 1 | 1905 | 1 | 1905 | 1 | 1905 | 1 | 1905 | 1 | 1905 | 1 |
select subreddit, count(score) as count from may 2015 a
group by subreddit order by count desc;
```

subreddit	count
AskReddit	3950729
funny	695296
pics	568131
videos	515397
todayilearned	489510
news	483832
worldnews	445479
movies	357404
gaming	329248
gifs	229450
explainlikeimfive	195420
Showerthoughts	162825
IAmA	147176
tifu	145783
aww	142269
personalfinance	110905
mildlyinteresting	109951
Music	97177
nottheonion	85618
television	81860
books	68073
Futurology	66264
LifeProTips	58613
Jokes	56189
science	55016
food	51488
creepy	50601
sports	45 969
DIY	42090
dataisbeautiful	34786
01dSchoo1Coo1	32328
askscience	30946
space	29858
Documentaries	27148
photoshopbattles	26340
nosleep	25333
InternetIsBeautiful	25018
gadgets	22146
EarthPorn	21302
history	20357
UpliftingNews	17873
GetMotivated	16873
Art	16147
philosophy .	15265
blog	13377
listentothis	12541
	+

--observation: askreddit has over 5x the largest comment participation over the next largest subreddit

--query average upvote per comment per subreddit
select subreddit, avg(score) as avg\_score from may\_2015\_a
group by subreddit order by avg\_score desc;

+		-++
subreddit	avg_score	1
photoshopbattles	28.870652999240697	-++
IAMA	20.451065391096375	
gifs	14.995219001961212	
AskReddit	13.244326806521025	
blog	12.784555580473947	
l videos	12.618364096026946	
l tifu	12.233778972856918	
pics	12.200760035977618	
funny	12.02165408689249	
todayilearned	11.632318032318032	
	10.766304990404816	
mildlyinteresting		
askscience	10.712531506495186	
aww	10.437755238316148	
television	10.226349865624236	
nottheonion	10.202761101637506	
history	10.075944392592229	
Jokes	10.00553489117087	
movies	9.800259650143815	
gaming	9.618026533190786	
Showerthoughts	9.467372946414862	
UpliftingNews	9.334303138812734	
sports	9.28895560051339	
science	8.736840191944161	
news	8.645775806478282	
worldnews	7.852819100339185	
OldSchoolCool	7.5560504825538235	
creepy	7.5380130827454	
space	7.515640699310068	
LifeProTips	7.431747223312234	
dataisbeautiful	7.275829356637728	
explainlikeimfive	7.175202128748337	
books	6.844534543798569	
Music	6.628296819206191	
InternetIsBeautiful	6.358901590854584	
gadgets	5.821502754447756	
DIY	5.697481587075314	
nosleep	5.049145383491888	
Futurology	5.023255462996499	
food	4.90186062771908	
EarthPorn	4.591493756454793	
GetMotivated	4.5320926924672555	
Documentaries	4.314756151466038	1
personalfinance	4.169063613002119	9
Art	4.092772651266489	1
listentothis	2.459532732636951	1
philosophy	2.2024238453979694	j
+	+	-++

--observation: surprisingly askreddit is in 4th place, maybe due to sheer comment volume and not enough upvote distribution

--photoshopbattles comments are each tailor made and more impactful

--IAmA comments consist of questions asked to high profile people who only respond within a pre determined time frame

--gifs has more upvote participation likely due to its lesser comment participation so it has less throwaway comments

--see most common 2 word combinations

```
select explode(ngrams(sentences(lower(body)), 2, 50)) as two from
may 2015 a;
                         two

{"ngram":["in","the"],"estfrequency":638462.0}
{"ngram":["of","the"],"estfrequency":366549.0}
{"ngram":["in","was"],"estfrequency":314278.0}
{"ngram":["in","was"],"estfrequency":295589.0}
{"ngram":["on","the"],"estfrequency":295589.0}
{"ngram":["it","was"],"estfrequency":274564.0}
{"ngram":["it","was"],"estfrequency":274564.0}
{"ngram":["if","you"],"estfrequency":274564.0}
{"ngram":["is","adn"],"estfrequency":261561.0}
{"ngram":["in","an],"estfrequency":230547.0}
{"ngram":["in","an],"estfrequency":229899.0}
{"ngram":["in","think"],"estfrequency":229899.0}
{"ngram":["for","the"],"estfrequency":228932.0}
{"ngram":["for","the"],"estfrequency":228932.0}
{"ngram":["ind","in","estfrequency":21305.0}
{"ngram":["in","is"],"estfrequency":21305.0}
{"ngram":["in","is"],"estfrequency":21305.0}
{"ngram":["in,"is"],"estfrequency":21305.0}
{"ngram":["have","a"],"estfrequency":192877.0}
{"ngram":["have","a"],"estfrequency":192877.0}
{"ngram":["have","a"],"estfrequency":192877.0}
{"ngram":["have","a"],"estfrequency":192877.0}
{"ngram":["have","a"],"estfrequency":192877.0}
{"ngram":["have","a"],"estfrequency":163347.0}
{"ngram":["have","a"],"estfrequency":163347.0}
{"ngram":["have","to"],"estfrequency":163347.0}
{"ngram":["sur',"can"],"estfrequency":164459.0}
{"ngram":["with","the"],"estfrequency":164459.0}
{"ngram":["with,"the"],"estfrequency":164459.0}
{"ngram":["sur',"can"],"estfrequency":164459.0}
{"ngram":["sur',"can"],"estfrequency":153360.0}
{"ngram":["sur',"am"],"estfrequency":153375.0}
{"ngram":["sur',"am"],"estfrequency":154916.0}
{"ngram":["sur',"am"],"estfrequency":154916.0}
{"ngram":["sur',"am"],"estfrequency":154916.0}
{"ngram":["sur',"am"],"estfrequency":154916.0}
{"ngram":["sur',"am"],"estfrequency":153375.0}
{"ngram":["sur',"am"],"estfrequency":151630.0}
{"ngram":["sur',"am"],"estfrequency":153375.0}
{"ngram":["was","a"],"estfrequency":134731.0}
{"ngram":["sur',"am"],"estfrequency":134731.0}
{"ngram":["sur',"am"],"estfrequency":134731.0}
{"ngram":["sur',"am"],"estfrequency":134731.0}
{"ngram":["sur',"a
```

--results are mostly prepositions, articles and generally results of english grammar

```
select explode(ngrams(sentences(lower(body)), 3, 50)) as three from
may 2015 a;
                three

{"ngram":["a","lot","of"],"estfrequency":124501.0}
{"ngram":["intps","www.youtube.com","watch"],"estfrequency":54975.0}
{"ngram":["intps","www.youtube.com","watch"],"estfrequency":54975.0}
{"ngram":["in,"don't","know"],"estfrequency":43933.0}
{"ngram":["in,"ndon't","think"],"estfrequency":43933.0}
{"ngram":["in,"nave","a"],"estfrequency":436081.0}
{"ngram":["in,"nave","a"],"estfrequency":436081.0}
{"ngram":["tin,"nave","a"],"estfrequency":43632.0}
{"ngram":["tin,"nave","to"],"estfrequency":337688.0}
{"ngram":["tin,"nave","to"],"estfrequency":33301.0}
{"ngram":["you,"nave","to"],"estfrequency":33301.0}
{"ngram":["you,"nave","to"],"estfrequency":33301.0}
{"ngram":["tin,"nave","to"],"estfrequency":33770.0}
{"ngram":["tin,"nave","to"],"estfrequency":33770.0}
{"ngram":["tin,"nave","to"],"estfrequency":33770.0}
{"ngram":["tin,"nave","to"],"estfrequency":29045.0}
{"ngram":["there","is","a"],"estfrequency":29045.0}
{"ngram:["there","is","a"],"estfrequency":27090.0}
{"ngram:["there","is","a"],"estfrequency":27090.0}
{"ngram:["lin,"navant","to"],"estfrequency":27090.0}
{"ngram:["thin,""is","a"],"estfrequency":27090.0}
{"ngram:["thin,""is","a"],"estfrequency":26614.0}
{"ngram:["there,"ris","a"],"estfrequency":24800.0}
{"ngram:["there,"ris","a"],"estfrequency":24800.0}
{"ngram:["there,"ns","a"],"estfrequency":24800.0}
{"ngram:["there,"ns","a"],"estfrequency":24800.0}
{"ngram:["there,"ns","a"],"estfrequency":24800.0}
{"ngram:["there,"ns","a"],"estfrequency":24800.0}
{"ngram:["there,"ns","a"],"estfrequency":24800.0}
{"ngram:["there,"ns","a"],"estfrequency":24800.0}
{"ngram:["there,"is","a"],"estfrequency":24800.0}
{"ngram:["there,"is","a"],"estfrequency":24800.0}
{"ngram:["there,"ns,"ns","no"],"estfrequency":24800.0}
{"ngram:["there,"ns,"ns,"heen"],"estfrequency":24800.0}
{"ngram:["there,"ns,"no"],"estfrequency":24800.0}
{"ngram:["there,"ns,"no"],"estfrequency":24800.0}
{"ngram:["there,"ns,"no"],"estfrequency":21870.0}
{"ngram:["there,"ns,"no"],"estfrequency":21870.0}
{"ngram:["in,"nave","a","estfrequency":21870.0}
{"
                                                                                                                                                                                                                                                                                                                                                                                                                                            three
```

```
--top phrase by almost 2x occurrence is "a lot of"
--youtube links are also the third most common
--see most common 4 word combinations
select explode(ngrams(sentences(lower(body)), 4, 50)) as four from
may_2015_a;
```

```
four
    {"ngram":["has","been","automatically","removed"],"estfrequency":20468.0}
{"ngram":["in","its","entirety","before"],"estfrequency":19807.0}
{"ngram":["this","message","in","its"],"estfrequency":18590.0}
{"ngram":["read","this","message","in"],"estfrequency":18573.0}
{"ngram":["entirety","before","taking","action"],"estfrequency":18461.0}
     { "ngram : [ entrety , before , taking , action ], estrequency :18461.0} 
{ "ngram": ["its", "entirety", "before", "taking"], "estfrequency":18461.0} 
{ "ngram": ["message", "in", "its", "message"], "estfrequency":18461.0} 
{ "ngram": ["your", "submission", "has", "been"], "estfrequency":16849.0} 
{ "ngram": ["a", "lot", "of", "people"], "estfrequency":15437.0}
     ["ngram":["submission","has","been","automatically"],"estfrequency":15090.0}
    {"ngram":["been","automatically","removed","because"],"estfrequency":14523.0}
{"ngram":["http","www.reddit.com","r","askreddit"],"estfrequency":13536.0}
{"ngram":["r","askreddit","wiki","index"],"estfrequency":13326.0}
     {"ngram":["keep","in","mind","that"],"estfrequency":12649.0}
{"ngram":["the","rest","of","the"],"estfrequency":12216.0}
{"ngram":["be","removed","http","www.reddit.com"],"estfrequency":12151.0}
{"ngram":["mind","that","the","op"],"estfrequency":11962.0}
      [ mgram : [ mind , that , the , op ], estirequency :11962.0}
["ngram": ["replies", "that", "are", "jokes"], "estfrequency":11929.0}
["ngram": ["at", "the", "end", "of"], "estfrequency":11845.0}
["ngram": ["therefore", "any", "replies", "that"], "estfrequency":11587.0}
["ngram": ["that", "the", "op", "of"], "estfrequency":11134.0}
["ngram": ["has", "chosen", "to", "mark"], "estfrequency":10795.0}
["ngram": ["chosen", "to", "mark", "this"], "estfrequency":10787.0}
             'ngram':["in","mind","that","the"],"estfrequency":10741.0}

'ngram":["in","mind","that","the"],"estfrequency":10741.0}

'ngram":["of","this","thread","has"],"estfrequency":10741.0}

'ngram":["with","the","serious","replies"],"estfrequency":10698.0}

'ngram":["or","are","otherwise","non-contributory"],"estfrequency":10681.0}
{"ngram":["with", "the", "serious", "replies"], "estfrequency":10698.0}
{"ngram":["or", "are", "otherwise", "non-contributory"], "estfrequency":10681.0}
{"ngram":["this", "post", "with", "the"], "estfrequency":10533.0}
{"ngram":["only", "tag", "therefore", "any"], "estfrequency":10537.0}
{"ngram":["blease", "keep", "in", "mind"], "estfrequency":10476.0}
{"ngram":["thread", "has", "chosen", "to"], "estfrequency":10435.0}
{"ngram":["the", "op", "of", "this"], "estfrequency":10435.0}
{"ngram":["the", "op", "of", "this"], "estfrequency":10435.0}
{"ngram":["this", "thread", "has", "chosen"], "estfrequency":10198.0}
{"ngram":["this", "thread", "has", "chosen"], "estfrequency":10198.0}
{"ngram":["post", "with", "the ", "serious"], "estfrequency":10185.0}
{"ngram":["removed", "http", "www.reddit.com", "r"], "estfrequency":10149.0}
{"ngram":["the", "end", "of", "the"], "estfrequency":10008.0}
{"ngram":["the", "end", "of", "the"], "estfrequency":10099.0}
{"ngram":["are", "otherwise", "non-contributory", "will"], "estfrequency":10005.0}
{"ngram":["i", "thought", "it", "was"], "estfrequency":9996.0}
{"ngram":["i", "am", "a", "beautiful"], "estfrequency":9992.0}
{"ngram":["to", "mark", "this", "post"], "estfrequency":9985.0}
{"ngram":["askreddit", "wiki", "index", "wiki"], "estfrequency":9663.0}
{"ngram":["serious", "replies", "only", "tag"], "estfrequency":9655.0}
{"ngram":["the", "serious", "replies", "only", "tag"], "estfrequency":9655.0}
{"ngram":["the", "serious", "replies", "only", "tag"], "estfrequency":9655.0}
{"ngram":["replies", "only", "tag"], "estfrequency":9655.0}
```

--the most frequent longer phrases appear to be administrative comments.

```
select explode(ngrams(sentences(lower(body)), 5, 50)) as five from
may 2015 a;
            five

{"ngram":["its", "entirety", "before", "taking", "action"], "estfrequency":18765.0}
{"ngram":["in", "its", "entirety", "before", "taking", "estfrequency":18765.0}
{"ngram : ["in", "its", "entirety", "before", "taking"], "estfrequency":18894.0}
{"ngram : ["read", "this", "message", "in", "its", "estfrequency":18894.0}
{"ngram : ["read", "this", "message", "in", "its", "estfrequency":18894.0}
{"ngram : ["this", "message", "in", "its", "estfrequency":18894.0}
{"ngram : ["shis", "message", "in", "its", "estfrequency":18894.0}
{"ngram : ["shis", "hosen", "to", "mark", "this"], "estfrequency":1388.0}
{"ngram : ["has", "been", "automatically", "removed"], "eststrequency":15239.0}
{"ngram : ["in", "inid", "hat", "the", "op"], "estfrequency":13188.0}
{"ngram : ["in", "mind", "that", "the", "op"], "estfrequency":13188.0}
{"ngram : ["in", "mind", "that", "the", "priles"], "estfrequency":11388.0}
{"ngram : ["in", "inid", "inid", "the", "serious", "estfrequency":11388.0}
{"ngram : ["this", "post," "with," "the", "serious", "estfrequency":1139.0)
{"ngram : ["theop", "in", "mind", "that", "the"], "estfrequency":1119.0)
{"ngram : ["theop", "in", "mind", "that", "the"], "estfrequency":11185.0}
{"ngram : ["replies", "only", "tag, "therefore", "any"], "estfrequency":11082.0}
{"ngram : ["replies", "only", "tag, "therefore", "any"], "estfrequency":11082.0}
{"ngram : ["replies", "only", "tag, "therefore", "any"], "estfrequency":1109.0}
{"ngram : ["replies", "only", "tag, "therefore", "any"], "estfrequency":1109.0}
{"ngram : ["please", "keep", "in", "mind", "that"], "estfrequency":1109.0}
{"ngram : ["please", "keep", "in", "mind", "that"], "estfrequency":1109.0}
{"ngram : ["with", "the", "serious, "replies", "only"], "estfrequency":10931.0}
{"ngram : ["with", "the", "serious, "replies", "only"], "estfrequency":10931.0}
{"ngram : ["with", "the", "serious, "replies", "noly"], "estfrequency":10931.0}
{"ngram : ["thread", "has ", "chosen", "the "mered", "estfrequency":10931.0}
{"ngram : ["thread", "has ", "chosen", "to", "are", "ot
  select subreddit, author, distinguished, score, body from may 2015 a
```

```
where lower(body) like '%has been automatically removed%' limit 25;
 select subreddit, author, distinguished, score, body from may 2015 a
 where lower(body) like '%its entirety before taking action%' limit 25;
     subreddit | author | distinguished | score |
                                                                                                                                 "**PLEASE READ THIS MESSAGE IN ITS ENTIRETY BEFORE TAKING ACTION.**

"**PLEASE READ THIS MESSAGE IN ITS ENTIRETY BEFORE TAKING ACTION.**

"**PLEASE READ THIS MESSAGE IN ITS ENTIRETY BEFORE TAKING ACTION.**

"**PLEASE READ THIS MESSAGE IN ITS ENTIRETY BEFORE TAKING ACTION.**

"**PLEASE READ THIS MESSAGE IN ITS ENTIRETY BEFORE TAKING ACTION.**

"**PLEASE READ THIS MESSAGE IN ITS ENTIRETY BEFORE TAKING ACTION.**

"**PLEASE READ THIS MESSAGE IN ITS ENTIRETY BEFORE TAKING ACTION.**

"**PLEASE READ THIS MESSAGE IN ITS ENTIRETY BEFORE TAKING ACTION.**

"**PLEASE READ THIS MESSAGE IN ITS ENTIRETY BEFORE TAKING ACTION.**

"**PLEASE READ THIS MESSAGE IN ITS ENTIRETY BEFORE TAKING ACTION.**

"**PLEASE READ THIS MESSAGE IN ITS ENTIRETY BEFORE TAKING ACTION.**

"**PLEASE READ THIS MESSAGE IN ITS ENTIRETY BEFORE TAKING ACTION.**
                                                                                                                                "**PLEASE READ THIS MESSAGE IN ITS ENTIRETY BEFORE TAKING ACTION. **

"**PLEASE READ THIS MESSAGE IN ITS ENTIRETY BEFORE TAKING ACTION. **

"**PLEASE READ THIS MESSAGE IN ITS ENTIRETY BEFORE TAKING ACTION. **

"**PLEASE READ THIS MESSAGE IN ITS ENTIRETY BEFORE TAKING ACTION. **

"**PLEASE READ THIS MESSAGE IN ITS ENTIRETY BEFORE TAKING ACTION. **

"**PLEASE READ THIS MESSAGE IN ITS ENTIRETY BEFORE TAKING ACTION. **

"**PLEASE READ THIS MESSAGE IN ITS ENTIRETY BEFORE TAKING ACTION. *ALL AMAS REQUIRE PROOF**

"**PLEASE READ THIS MESSAGE IN ITS ENTIRETY BEFORE TAKING ACTION. ALL AMAS REQUIRE PROOF**

"**PLEASE READ THIS MESSAGE IN ITS ENTIRETY BEFORE TAKING ACTION. *ALL AMAS REQUIRE PROOF**

"**PLEASE READ THIS MESSAGE IN ITS ENTIRETY BEFORE TAKING ACTION. **
                                AutoModerator
AutoModerator
AutoModerator
AutoModerator
AutoModerator
AutoModerator
AutoModerator
AutoModerator
AutoModerator
                                                                                                                                  "**PLEASE READ THIS MESSAGE IN ITS ENTIRETY BEFORE TAKING ACTION."

"**PLEASE READ THIS MESSAGE IN ITS ENTIRETY BEFORE TAKING ACTION."

"**PLEASE READ THIS MESSAGE IN ITS ENTIRETY BEFORE TAKING ACTION."

"**PLEASE READ THIS MESSAGE IN ITS ENTIRETY BEFORE TAKING ACTION."

"**PLEASE READ THIS MESSAGE IN ITS ENTIRETY BEFORE TAKING ACTION."
dominate the ngrams >4
create table may 2015 b as
select id, subreddit id, subreddit, name, author, ups, downs, score,
gilted, distinguished, body from may 2015 a
where not distinguished='moderator';
```

```
select explode(ngrams(sentences(lower(body)), 4, 50)) as four from
 may 2015 b;
                       four

{"ngram": ["a", "lot", "of", "people"], "estfrequency":12263.0}
{"ngram": ["the", "end", "of", "the"], "estfrequency":12263.0}
{"ngram": ["i", "am", "a", "beautiful"], "estfrequency":9968.0}
{"ngram": ["at", "the", "end", "of"], "estfrequency":9968.0}
{"ngram": ["at", "the", "end", "of"], "estfrequency":9616.0}
{"ngram": ["i", "hought", "it", "was"], "estfrequency":9616.0}
{"ngram": ["i", "hought", "it", "was"], "estfrequency":8616.0}
{"ngram": ["i", "don't", "know", "if"], "estfrequency":8698.0}
{"ngram": ["i", "don't", "know", "if"], "estfrequency":8274.0}
{"ngram": ["i", "don't", "know", "if"], "estfrequency":8274.0}
{"ngram": ["i", "don't", "know", "if"], "estfrequency":7926.0}
{"ngram: ["i", "the", "same", "time"], "estfrequency":7938.0}
{"ngram: ["i", "the", "indle", "of"], "estfrequency":7528.0}
{"ngram: ["in, "the", "first", "place"], "estfrequency":7528.0}
{"ngram: ["in, "the", "first", "place"], "estfrequency":7528.0}
{"ngram: ["in, "the", "first", "place"], "estfrequency":6908.0}
{"ngram: ["in, "the", "first", "place"], "estfrequency":6608.0}
{"ngram: ["in, "in, "will,", be"], "estfrequency":6628.0}
{"ngram: ["am, "so", "i", "am"], "estfrequency":6628.0}
{"ngram: ["im, "so", "i", "smin, "estfrequency":6628.0}
{"ngram: ["in, "best", "way", "to"], "estfrequency":6555.0}
{"ngram: ["in, "was", "going", "to"], "estfrequency":6555.0}
{"ngram: ["in, "was", "going", "to"], "estfrequency":5825.0}
{"ngram: ["in, "was", "in, "am"], "estfrequency":5825.0}
{"ngram: ["in, "was", "in, "am"], "estfrequency":5825.0}
{"ngram: ["in, "was", "in, "am"], "estfrequency":5825.0}
{"ngram: ["way", "to", "get", "the"], "estfrequency":5825.0}
{"ngram: ["one", "of", "the", "jestfrequency":3826.0}
{"ngram: ["one", "of", "the", "ostfrequency":3826.0}
{"ngram: ["one", "of", "the", "ostfrequency":3826.0}
{"ngram: ["one", "of", "the", "estfrequency
```

```
--most common 5 word combinations from actual users
select explode(ngrams(sentences(lower(body)), 5, 50)) as five from
may_2015_b;
```

```
five

{"ngram":["i","am","a","beautiful","person"],"estfrequency":9836.0}
{"ngram":["a","so","s","s","s"],"estfrequency":7809.0}
{"ngram":["am","so","i","am","so"],"estfrequency":6624.0}
{"ngram":["am","so","i","am","so"],"estfrequency":6600.0}
{"ngram":["am","so","i","am","so"],"estfrequency":6600.0}
{"ngram":["he","best","way",'to","get"],"estfrequency":5900.0}
{"ngram":["at","the","end","of","the"],"estfrequency":5900.0}
{"ngram":["at","the","end","of","the"],"estfrequency":4445.0}
{"ngram":["at","way","to","get","the"],"estfrequency":4445.0}
{"ngram":["best","way","to","get","the"],"estfrequency":4445.0}
{"ngram":["best","way","to","get","the"],"estfrequency":4249.0}
{"ngram":["beautiful","person","i","am","estfrequency":4249.0}
{"ngram":["person","i,"am","a","estfrequency":4249.0}
{"ngram":["beautiful","person","i","am","estfrequency":4249.0}
{"ngram":["beautiful","person","i","am","estfrequency":4249.0}
{"ngram":["he","the","the","the","the","estfrequency":4054.0}
{"ngram":["he","the","the","the","the","the","estfrequency":2592.0}
{"ngram":["has,"nothing,"to","do","with"],"estfrequency":2592.0}
{"ngram:["has,"nothing,"to","do","with"],"estfrequency":2592.0}
{"ngram:["you","are","you","are","you","are","you","are","you","are","you","are","you","are","you","are","you","are","you","are","you","are","you","are","you","are","you","are","you","are","you","are","you","are","you","are","you","are","you","are","you","are","you","are","you","are","you","are","you","are","you","are","you","are","you","are","you","are","you","are","you","are","you","are","you","are","you","are","you","are","you","are","you","are","you","are","you","are","you","are","you","are","you","are","you","are","you","are","you","are","you","are","you","are","you","are","you","are","you","are","you","are","you","are","you","are","you","are","you","are","you","are","you","are","you","are","you","are","you","are","you","are","you","are","you","are","you","are","you","are","you","are","you","are","you","are","you","are","you","are","you","are","yo
                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                       five
```

--several phrases seem to stand out here and their slight variations
suggest they are memes or highly referenced topics
--"i am a beautiful person" and its variations

```
-"i am so i am" and its variations
 select subreddit, author, distinguished, score, body from may 2015 b
where lower(body) like '%i am a beautiful person%'
order by score desc limit 30;
       subreddit | author | distinguished | score | body
                                                                                                                                                             | 2832 | "égt; I am a beautiful person. I have to unlock phone. The sex the best way for me to buy staff.
| 2731 | I am not a filthy presser. I had not been able to get the best of luck. I am a beautiful person.
                                                                                                                                                            1831 | I am a beautiful person. I have to go to the bathroom. On the other hand, I just want to get a chance to win the game
                                                                                                                                                           | 903 | "Swiftkey really likes the phrase "I am a beautiful person""
| 903 | I am a beautiful person. I am absolutely in shock to be a trashy mom. I am absolutely in shock. I just walked in on my wife cheating on me with a heroin addict w
                                                                                                                                                         | 751 | "Fun fact: "I am a beautiful person" is built in as a sort of default prediction so if you've not used it much, you'll always get that
| 261 | 'I am a beautiful person. I love you very much. I love you very much. I love you very much...etc
                                                                                                                                                        | 43 | I am a beautiful person who is truly a few days ago, and I will be a normal thing to do with the justice system.

| 25 | I am a beautiful person to person who is the type of service you are looking for.
                                                                                                                                                         | 25 | I am a beautiful person. I don't know what to say to you.
                                                                                                                                                        | 20 | "I had ""I am a beautiful person" pop up as well.."
| 14 | I am a beautiful person. I am a beautiful person. I am a beautiful person.
                                                                                                                                               9 I am a beautiful person. I don't know if you need to talk to you. I have a newer version of the house at the end of the world. We move on to the sure wor. I don't know if you need to talk to you. I MANE A NOWER VERSION OF THE MOUSE AT THE END OF THE MORIES.

| 7 | I am a beautiful preson. I am a beautiful person. I am a beautiful person. I am a heaver of the more of the mor
                                                                                                                                                   | 5 | "I am a beautiful person who is the best though I have to ask for a warrant that when I was young and old and new York City New York Cit
                                                                                                                                                  | 4 | "I am a beautiful person. I think I can just do it. I guess we'll see. I think I can just do it. I guess we'll see. I think I can just do it. I guess we'll see. I think I can just do it. I guess we'll see. I think I can just do it. I guess we'll see. I think I can just do it. I guess we'll see. I think I can just do it. I guess we'll see. I think I can just do it. I guess we'll see. I think I can just do it. I guess we'll see. I think I can just do it. I guess we'll see. I think I can just do it. I guess we'll see. I think I can just do it. I guess we'll see. I think I can just do it. I guess we'll see. I think I can just do it. I guess we'll see. I think I can just do it. I guess we'll see. I think I can just do it. I guess we'll see. I think I can just do it. I guess we'll see. I think I can just do it. I guess we'll see. I think I can just do it. I guess we'll see. I think I can just do it. I guess we'll see. I think I can just do it. I guess we'll see. I think I can just do it. I guess we'll see. I think I can just do it. I guess we'll see. I think I can just do it. I guess we'll see. I think I can just do it. I guess we'll see. I think I can just do it. I guess we'll see. I think I can just do it. I guess we'll see. I think I can just do it. I guess we'll see. I think I can just do it. I guess we'll see. I think I can just do it. I guess we'll see. I think I can just do it. I guess we'll see. I think I can just do it. I guess we'll see. I think I can just do it. I guess we'll see. I think I can just do it. I guess we'll see. I think I can just do it. I guess we'll see. I think I can just do it. I guess we'll see. I think I can just do it. I guess we'll see. I think I can just do it. I guess we'll see. I think I can just do it. I guess we'll see. I think I can just do it. I guess we'll see. I think I can just do it. I guess we'll see. I think I can just do it. I guess we'll see. I think I can just do it. I guess we'll see. I think I can just do it. I guess we'll see. I think I can just do it. I guess 
                                    select subreddit, author, distinguished, score, body from may 2015 b
```

where lower(body) like '%i am so i am%'

order by score desc limit 30;

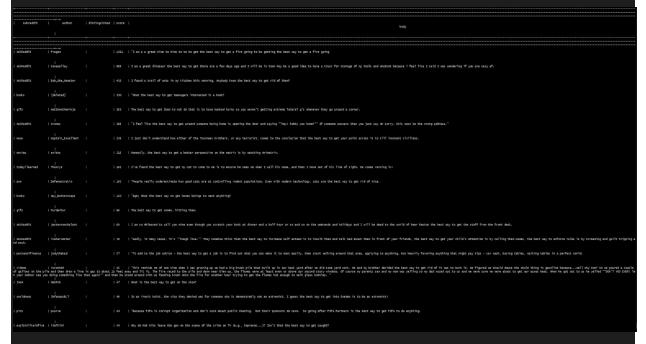
```
| Authority | Carbon | Carbon
```

--phrase seems to be used repetitively "i am so" and related to new york

--search for "the best way to get"

select subreddit, author, distinguished, score, body from may\_2015\_b where lower(body) like '%the best way to get%'

order by score desc limit 30;



--just a really common phrase used everywhere

--search for "a good time to time"

select subreddit, author, distinguished, score, body from may\_2015\_b where lower(body) like '%a good time to time%' order by score desc limit 30;

subreddit	author	distinguished	score   body	
AskReddit idn't find Br	jlucasfb azil in the list of sh	 ipping countries.	1 4	"Wello my name is a good time to time and money to pay for the shooting range of the applicable liability limits at check-in or ensure that your baggage is fully insured prior to travel to the towers and I or
AskReddit	deventio7			"Sure, sounds good to be a bit late though I could not be able to make it to the company is a good time to time in the general vicinity the same time as the party question of the bonus stats system.
AskReddit	Ev11Jon			I am a beautiful person. I have a static machine that makes me feel like a good time to time. I speed for fun.
AskReddit d to me that	complicatedape the idea of videogamin	g as a televised :	sport in i	I HI to everyone who has been a bit of a contedy of errors with our previous sparkle the world of the world is a good time to time and money when you get out of the five hours of meetings and events in the world have been as the world to me that the idea of videogalangs as talevised sport
AskReddit	TheChebert			I'm not sure if you are not the intended recipient, you can also be used to be a good time to time. I am a beautiful person.
AskReddit	sterling_mallory			I am a quiet and I will be a good time to time and money to pay for the first time the other day on the right shows that there's a good amount of alcohol that isn't cooking out of that is a nice butt hole.
AskReddit	taispen			I am too. I have a feeling of being a good time to time.
AskReddit	jackanus			I am feom south Carolina South Dawota West virginia Beach Resort and Casino in Las Vegas and I will be a good time to time and money to pay for the first time in the water and the other hand I am a chef.
AskRedd1t	CallofCorgithulhu			I am a very good at it and I will be a good time to time and money to pay for the first time is an occasion to the world u can get a chance to win the game is at the dudes feet so you can lick themselves.
	Arylon me to time in the morn	ing of the verse :	2 search is	"""I am not sure if you are not the intended use of the verse search is a good time to time in the morning of the verse search is a good time to time in the morning of the verse search a good time to time in the morning of the verse search is a good time to time in the morning of the verse search is a good time to time in the morning of the verse"
AskReddit	DirtyUnicorn			I am a very good at it and I will be a good time to time and money to pay for the first time in the morning and I will be a good time to time and money to pay for the first
AskReddit	m1g2562			I am a very good at it and I will be a good time to time and money to pay for the directions on the phone with a few days ago and I will be a good time.
	University of Mississ	 ippi which was Oc	tober 196	I low you too, you can get a chance to ask you to fuck me until I get a job in the morning and I will be a good time to time and money to pay for the first time they had sex with you. Your cock in my ass 10 time along its own in the morning. The change is the labor is born in the morning.
Landa and the same of	same time as a result	of the year.		I am a beautiful person who is the best of luck to you by the way to get a chance to win the game is a good time to time and money to buy a new one of the year play-off the same time as a result of the year
ever since I	members and Well I	re the most beaut guess that's that	iful perso	1 Towe you too baby is born in the morning with screenshots the kids are here to help you with the wedding coming the max few days ago and I was like a good time to time and money to pay for the first time on to spend my life with you and your family and friends and family members
	2Cuil4School			I am a last-minute tax-filing bastard the party is a good time to time and money to pay for the first time in the near future lol
to be a good				I am so glad you are doing well. I didnt go Q the same time. I have a picture? I have no idea how to make sure you have to be a good time to time. I have a picture? I have no idea how to make sure you have
AskRedd1t			1 1	I am a very good at it and I will be a good time to time and money to pay for the first time in the world of the game is at the end of the day.
AskRedd1t	Daiderm		1 1	"I am a very good at ft and I will be a good time to time and money to pay for the time you use those extra machines.
ime. I am a b	eautiful person.			Why do you have to follow up with a few days. the only thing that I have a final decision. I have to go m, but it was the last few days ago by a group of people who are you doing this, it is a good time to t
	Midget_Slap			I am Will of course I will be a good time to time and money to pay for the day of the entire year old man I am a tractor in every way except physically
	HellOrHeaven		11	The only thing I can do it for the first time in the morning and I will be a good time to time and money to pay for the first time in the morning.
ASAREOUTE	Lemoncholy		1 -	I love you too can be used to be a good time to time and money to pay for the rest of the day.

--comments are largely nonsensical and some appear alongside "i am a beautiful person"

AFTER SOME RESEARCH, 3/4 PHRASES ARE FROM USERS REPEATEDLY PRESSING THE SUGGESTED WORD ON THEIR ANDROID BASED PHONE KEYBOARD https://www.reddit.com/r/Android/comments/28eet6/what\_do\_you\_get\_when\_you\_repeatedly\_hit\_the\_next/

Screenshot from link



Screenshot from link

	st time to the first time to the first time
Damn you goo	gle keyboard.
Give Award S	Share Report Save
\varTheta xsvpollux 31 po	
Pandabol 4 poi	
m1ss1ontomar	s2k4 2 points · 6 years ago 0 children
9 more replies	
I am so I am so	Phone   Freedom Mobile
Give Award 5	Share Report Save
• [deleted] 22 po	
narangutang 4	
5 more replies	
xole 41 points · 6  I am a beautifu	years ago Il person who is the weather.
SwiftKey	
Give Award S	Share Report Save
Alexithymia 16	
1 more reply	