

Introduction to sentiment analysis

Whenever we need to make a decision, we often seek out the opinions of others. In the past, we seek opinions from friends and family or companies would use surveys, focus groups, opinion polls, and consultants. Now, customer reviews on the Internet has risen exponentially over the last decade. It is an important resource for buying products or attending events. In these situations, we would like to see what others are saying about them. Also it is a significant aspect for companies making decisions about their products or services.

Sentiment analysis, the computational study of how opinions, attitudes and emotions are expressed in nature language, provides techniques for extracting some emotional words, subjective information from large datasets of customer reviews and summarizing it. It can thus be vital to service providers, product producers, moviemakers, allowing them to quickly assess how new products and features are being received.



This graph is an example for a sentiment word cloud

Objectives

Sentiment analysis involves cross-fields backgrounds. It is one type of data science. We can use it on several different fields such as movie industry, electronic product industry and some services industry.

- ✓ Learn the conception of sentiment analysis. Doing a good sentiment analysis need to prepare well in a lot fields. Dealing with the natural language is very hard. First is the text normalization. Then it should be the whitespace tokenizer. Third is the Bayes classifier.
- ✓ Gain real understanding about programming in R on data collection and implement sentiment analysis. By utilizing statistical software R to retrieve data from html or xml format.
- \checkmark The objective, while working on this independent study, will be to understand how to conduct research, do research survey and think through the problems and figuring out the path to the best result. Through this project I am also going to learn some Machine Learning techniques and practices.

ww.PosterPresentations.co

Sentiment Analysis for Ipad 2 Sida Ye

Prof. Aldous

<text><code-block><text><text><code-block><text></text></code-block></text></text></code-block></text>			Method	ls to	implem	nent	sentim	ent ana	lysis		
<text><code-block><text></text></code-block></text>		F so au A th fo	irst step is oftware, R, utomaticall mazon. I c he function ormat file.	colle to w ly retronly n will	ction data for called	From he fu mer r ut the custo	amazon. I nction, wh reviews and production mer review	used the s tich can d rates fro n ID in an ws and rate	tatistical m nazon and es in a csv		
<text></text>			<pre>2 - datacollecti 3 library(XM 4 url <- pas 5 url <- pas 6 url <- pas 7 doc <- htm 8 numreview 9 numpage <- 10 url <- gsu 11 url <- pas 12 reviews <- 13 - for(i in 1 14 doc<-htm 15 reviews[16 finalrev 17 } 18 rates<-lis 19 - for(i in 1 20 doc<-h 21 rates[22 } 23 numericrat</pre>	<pre>ion <- func ful ste0("http: ste0(url,"/ ste0(url,"1 nlParse(url <- as.nume - numreview ub("1","",u ste0(url,1: - list() l:numpage){ nlParse(url [i]<-list(x st() l:numpage){ ntmlParse(url st() l:numpage){ ntmlParse(url st()</pre>	<pre>tion(ProdID){ //www.amazon.com/prod /ref=cm_cr_pr_top_link ")) pric(gsub(",","",strsp %/%10 + 1 wrl) numpage) [] [i]) pathSApply(doc,'//div st(reviews) [] [i]) pathSApply(doc,'//space) [] []] [] [] [] [] [] [] [] [] [] [] []</pre>	uct-review _2?ie=UTF8 lit(xpaths [@class="r	vs/",ProdID) &&pageNumber=") &Apply(doc,'//b',xmlV reviewText"]',xmlValu	/alue)[1]," ") µe))) "]1//span[@cl			
<text><text><code-block></code-block></text></text>			24 + for (1 1n 25 numeri 26 } 27 dfreview<- 28 write.csv(<pre>l:numrevie crates[i]< data.frame dfreview.P</pre>	w){ -as.numeric(strsplit(:(ReviewID=c(1:length(ProdID)	unlist(rat finalrevie	es)," ")[[i]][1]) ws)),customer_review	vs=finalreviev			
<pre>the function, we get eight units.</pre>		Second is going to find the list of words, which contains both positive words and negative words. Then, the next step is to tokenization. Tokenizing (splitting a string into its desired constituent parts) is fundamental to all Neuro-linguistic programming tasks. It is very complex if you want to do it extremely accurate, because sometime it is unusually represented a single cluster of punctuation like :-(might already represent the whole opinion. Hence, I did in a sample way. I use the function strsplit in R in order to split every word with a space. For example, "I love inad. It is awesome," By applying									
 standoved -2 is also and over the standoved -2 is also also also also also also also als		tł	ne function	, we §	get eight ur	nits.					
 Then we are moving to the next step. I try to match the words from word list with words in the reviews. Here, I use a Bayesian classifier, which means one positive match counts one point, and one negative match counts negative one point. scores = laply(sentences, function(sentence, pos.words, neg.words) {		1 2 3 4 5 6 7 8 9 10	abandon abandoned abandons abducted abduction abductions abhor abhorred abhorrent abhors abilities	-2 -2 -2 -2 -2 -2 -2 -3 -3 -3 -3 -3 2	abandon abandoned abandons abducted abduction abductions abhor abhorred abhorrent abhors abilities	<pre>> str [[1]] [1] [7] [13] [19] [25] [31]</pre>	<pre>split(a[2]," "Exactly" "The" "is" "to" "connected." "if"</pre>	") "what" "cord" "crimped" "attach" "I" "at"	"shipper" "connecting" "so" "but" "won't" "all"		
<pre>scores = laply(sentences, function(sentence, pos.words, neg.words) { # clean up sentences with R's regex-driven global substitute, gsub(): sentence = gsub('[[:punct:]]', '', sentence) sentence = gsub('[cintr1:]]', '', sentence) # and convert to lower case: sentence = tolower(sentence) # split into words. str_split(sentence, '\\s+') # sometimes a list() is one level of hierarchy too much words = unlist(word.list) Split all the sentence into single, lower letter words. This is a preparation for the further sentiment analysis. # compare our words to the dictionaries of positive & negative terms pos.matches = match(words, pos.words) neg.matches = match(words, neg.words) # match() returns the position of the matched term or NA # we just wont a TRUE/FALSE: pos.matches = lis.na(pos.matches) # and conveniently enough, TRUE/FALSE will be treated as 1/0 by sum(): score = sum(pos.matches) - sum(neg.matches) Based on the method above, I can get a sentiment score finally. </pre>		T fr B p	hen we are om word l ayesian cla oint, and o	e mov ist wi assifie ne ne	ving to the f th words if er, which n gative mat	next s n the neans ch co	step. I try t reviews. H one positi unts negat	o match th Iere, I use ive match ive one po	ne words a counts one oint.		
 # split into words. str_split is in the stringr package word.list = str_split(sentence, '\\s+') # sometimes a list() is one level of hierarchy too much words = unlist(word.list) > Split all the sentence into single, lower letter words. This is a preparation for the further sentiment analysis. # compare our words to the dictionaries of positive & negative terms pos.matches = match(words, pos.words) neg.matches = match(words, neg.words) # match() returns the position of the matched term or NA # we just want a TRUE/FALSE: pos.matches = !is.na(pos.matches) neg.matches = !is.na(neg.matches) # and conveniently enough, TRUE/FALSE will be treated as 1/0 by sum(): score = sum(pos.matches) - sum(neg.matches) > Based on the method above, I can get a sentiment score finally. 	<pre>scores = laply(sentences, function(sentence, pos.words, neg.words) { # clean up sentences with R's regex-driven global substitute, gsub(): sentence = gsub('[[:punct:]]', '', sentence) sentence = gsub('[[:cntrl:]]', '', sentence) sentence = gsub('\\d+', '', sentence) # and convert to lower case: sentence = tolower(sentence)</pre>										
 Split all the sentence into single, lower letter words. This is a preparation for the further sentiment analysis. # compare our words to the dictionaries of positive & negative terms pos.matches = match(words, pos.words) neg.matches = match(words, neg.words) # match() returns the position of the matched term or NA # we just want a TRUE/FALSE: pos.matches = !is.na(pos.matches) neg.matches = !is.na(neg.matches) # and conveniently enough, TRUE/FALSE will be treated as 1/0 by sum(): score = sum(pos.matches) - sum(neg.matches) Based on the method above, I can get a sentiment score finally. 			<pre># split int word.list = # sometimes words</pre>	to words = str_sp s a list	<pre>s. str_split is lit(sentence, ' () is one level </pre>	in the \\s+') of hie	stringr package rarchy too much	e h			
<pre>pos.matches = match(words, pos.words) neg.matches = match(words, neg.words) # match() returns the position of the matched term or NA # we just want a TRUE/FALSE: pos.matches = !is.na(pos.matches) neg.matches = !is.na(neg.matches) # and conveniently enough, TRUE/FALSE will be treated as 1/0 by sum(): score = sum(pos.matches) - sum(neg.matches)</pre>		S p	plit all the reparation	sente for th	ence into sinte the further set	ngle, entim	lower lette ent analys	er words. 7 is.	This is a		
<pre># match() returns the position of the matched term or NA # we just want a TRUE/FALSE: pos.matches = !is.na(pos.matches) neg.matches = !is.na(neg.matches) # and conveniently enough, TRUE/FALSE will be treated as 1/0 by sum(): score = sum(pos.matches) - sum(neg.matches) > Based on the method above, I can get a sentiment score finally.</pre>		<pre>pos.matches = match(words, pos.words) neg.matches = match(words, neg.words)</pre>									
 # and conveniently enough, TRUE/FALSE will be treated as 1/0 by sum(): score = sum(pos.matches) - sum(neg.matches) > Based on the method above, I can get a sentiment score finally. 	<pre># match() returns the position of the matched term or NA # we just want a TRUE/FALSE: pos.matches = !is.na(pos.matches) neg.matches = !is.na(neg.matches)</pre>										
Based on the method above, I can get a sentiment score finally.			<pre># and conve score = sum</pre>	niently (pos.ma	enough, TRUE/F/ tches) - sum(neg	ALSE wil g.matche	ll be treated a	s 1/0 by sum()	:		
		B	ased on the	e met	hod above	, I cai	n get a sen	timent sco	ore finally.		

Department of Statistics

Results

This data frame contains the review ID, amazon star rate, sentiment score, number of positive words, number of negative words, length of reviews and the context of customer reviews

ReviewID	rates	sentiment_scores	numofpos	numofneg	length_of_words	customer_reviews
1	5	0	2	2	47	I was recently reminded how cool i
2	4	-1	0	1	34	Exactly what shipper described. Pr
3	5	-1	0	1	20	I bought this for my husband for h
4	3	-2	1	3	27	When I received the item, it was d
5	5	5	5	0	20	good product woks very well and it
6	5	4	4	0	53	Nextworth was very prompt and gave
7	5	4	5	1	46	Love love love! Came super fast, a
8	4	3	3	0	36	This unit was in great shape, it l
9	1	-1	0	1	57	this item its used!!!! it was alre
10	5	3	3	0	62	The IPad that I purchased was not
11	5	4	4	0	23	Love its flexibility and portabili
12	5	5	5	0	61	Bought used iPad 2 for family use.
13	5	A	1	1	30	Got the 64GR refurbished model for

Here is the rate distribution for ipad2





By comparing these two plots, we can see that the distribution of sentiment scores follows the distribution of amazon rates. As ipad2 has plenty of rates with 5 stars, the sentiment scores are also skew to left, which means that most of the sentiments in the reviews are positive. It implies that ipad2 is popular among the public.



From the results, we can conclude that ipad 2 has a very positive customer reviews. But I still want to know how popular it is. Thus, I decided to compare it with its competitors, Samsung Galaxy Tab 10 and HP Touchpad. I did the same steps for the other two products. Then, i got these pie plots below.

From the pie plots, we can see that the satisfaction of samsung > ipad > hp pad. By applying the sentiment analysis, though, it is not an accurate conclusion, but we could say that samsung tab 10 has the highest satisfaction rate among these three products. Due to the limitation of the poster, I just post part of my further idea about the regression with the data of ipad 2.



We can use the sentiment analysis to check whether the customers satisfy with companies' products, instead of reviewing the rates from Amazon or other online shopping websites. By implementing the sentiment analysis with only the context of customer reviews, we can figure out the feedbacks from customers. Thus, for further usage, we can analysis the data from tweet, which represents a more powerful resource.

Tutorial. guidance.



Conclusions





This is the word cloud for the customer reviews of ipad2 from Amazon References

• Potts, Christopher. "Sentiment Symposium Tutorial." Sentiment Symposium Acknowledgments

I want to appreciate Professor. Aldous for providing this opportunity to me to show my ideas about sentiment analysis and always providing encouraging