

Bag of Words Classification, Word Embeddings, Skip-Gram/CBOW Models Analysis

Mohammed Alrashidan

```
source("~/Users/mo/Desktop/Desktop/School/USF/Courses/Fall 2021/NLP/Functions/load_NLP_env.R")
path <- "~/Users/mo/Desktop/Desktop/School/USF/Courses/Fall 2021/NLP/Functions/"
load_NLP_env(path)

## [1] "functions loaded: "
```

```
library(Rcpp)
library(test2vec)
```

```
file_path <- "~/Users/mo/Desktop/Desktop/School/USF/Courses/Fall 2021/NLP/datasets/amazon_reviews.csv"
amazon <- read.csv(file_path)
amazon <- head(amazon, 20000)
```

Convert this to a binary classifier where overall scores of 4 and 5 are 'positive,' and other values are 'negative.'

```
library(dplyr)
amazonpos_neg <- ifelse(amazon$overall == 4 | amazon$overall == 5, 1, 0)
review_stat <- amazon %>% group_by(pos_neg) %>% tally()
review_stats_percent <- review_stat$pos_neg/review_stats$total*100
```

We can see that we have 81% 1s which is positive reviews compared to 19% of negative reviews

```
# cleaning the amazon reviews
text <- pre_process_corpus(amazon, "reviewText", replace_numbers = T, root_gm = "lemmatize")
amazonreview_preprocessed <- text
amazonreview_preprocessed[1]
```

```
## [1] "one thing book seem obvious original think however clarity author explain innovation happen remarkable al  
an program discuss mean human interaction kind situation tend inspire original clear think lead innovation thin  
g include people communicate certain situation outside normal pattern program identify ingredient make innovat  
on likely include people compel interact normally lead serendipity sometimes phenomenon will occur collaboration  
sometimes chance individual away home travel recommend book common sense truth apparent mastery subject author"
```

Construct a DTM that includes unigrams, bigrams, and trigrams.

```
# preparing the data for analysis.
# splitting the data to train and test
rand <- runif(nrow(amazon))
sets <- ifelse(rand < 0.9, "train", "test")
amazon$set <- sets

train <- amazon[amazon$set == "train",]

# constructing unigrams, bigrams, and trigrams.
dt_train <- tokenize(train$review_preprocessed, tokenizer = word_tokenizer, ids = train$id)
vocab <- create_vocabulary(dt_train, ngram = c(1,1))

lbound <- round(0.009 * nrow(train))
vocab <- vocab[vocab$doc_count > lbound,]
head(vocab)
```

```
## Number of docs: 189161
## 0 stopwords
## ngram_min = 1; ngram_max = 3
## Vocabulary:
##      term term_count doc_count
## 1:  fast_ship      1625      1625
## 2:  fit_perfect      1667      1667
## 3:   buy_now      1675      1675
## 4:   buy_sure      1684      1684
## 5:   satisfy      1685      1685
## 6:   guess      1763      1629
```

```
vectorizer <- vocab_vectorizer(vocab)
dtm_train <- create_dtm(dt_train, vectorizer)
dim(dtm_train)
```

```
## [1] 189161    425
```

```
test <- amazon[amazon$set == "test",]
it_test <- tokenize(train$review_preprocessed,
                    tokenizer = word_tokenizer, ids = test$id)
dtm_test <- create_dtm(it_test, vectorizer)
dim(dtm_test)
```

```
## [1] 18939    425
```

Regularization Method using Lasso

```
library(glmnet)
```

```
## Loaded glmnet 4.1-2
```

```
model_dtm <- cv.glmnet(x = dtm_train, y = train$pos_neg, type.measure = "auc",
                      family = "binomial", alpha = 1)
```

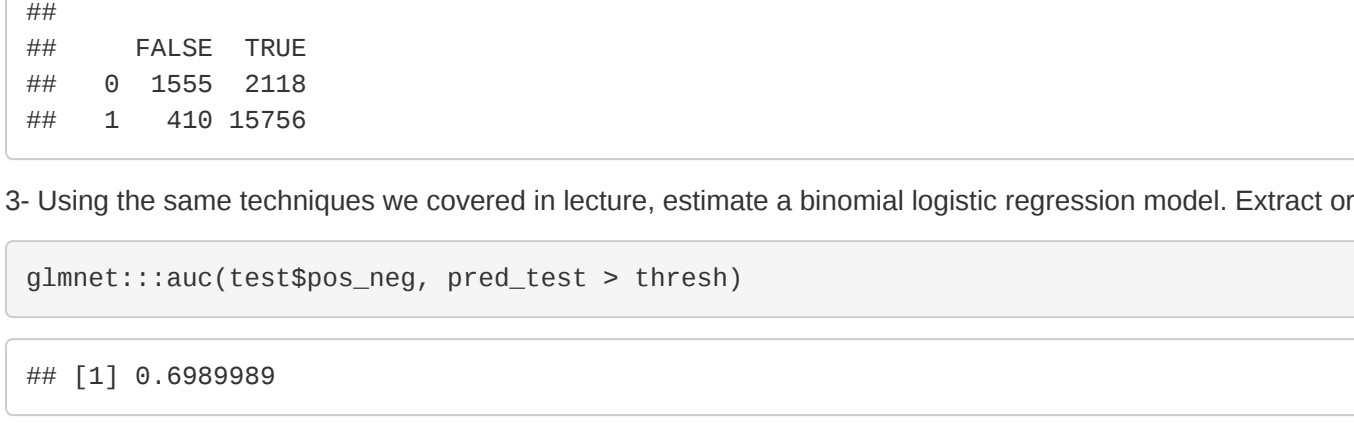
```
coefs <- coef(model_dtm, s = "lambda.min")
coefs <- data.frame(name = coefs$dimnames[1,1][coefs$ql + 1], coefficient = coefs$sk)
```

```
print((nrow(coefs)/ncol(dtm_train)))
```

```
## [1] 0.9411765
```

```
ggplot(coefs, aes(coefficient)) + geom_histogram(fill = "lightblue") + theme_classic()
```

```
## 'stat_bin()' using 'bins' = 30'. Pick better value with 'binwidth'.
```



```
pred_test <- predict(model_dtm, dtm_test, type= "response")[,1]
thresh <- 0.5
table(test$pos_neg, pred_test > thresh)
```

```
##
##      FALSE TRUE
## 0    1552 2118
## 1    419 15756
```

3- Using the same techniques we covered in lecture, estimate a binomial logistic regression model. Extract or calculate your model's AUC value.

```
glmnet::auc(test$pos_neg, pred_test > thresh)
```

```
## [1] 0.9089989
```

```
#model 2
model_dtm <- cv.glmnet(x = dtm_train,
                      y = train$pos_neg,
                      type.measure = "auc",
                      family = "binomial",
                      alpha = 1,
                      nfolds = 3)

coefs <- coef(model_dtm, s = "lambda.min")
coefs <- data.frame(name = coefs$dimnames[1,1][coefs$ql + 1], coefficient = coefs$sk)
```

```
print((nrow(coefs)/ncol(dtm_train)))
```

```
## [1] 0.9411765
```

```
pred_test <- predict(model_dtm, dtm_test, type= "response")[,1]
thresh <- 0.5
table(test$pos_neg, pred_test > thresh)
```

```
##
##      FALSE TRUE
## 0    1522 2161
## 1    361 15775
```

```
glmnet::auc(test$pos_neg, pred_test > thresh)
```

```
## [1] 0.693733
```

In this model we notice that decreasing the folds will decrease the AUC which means number of iteration should be higher

```
#model 3
model_dtm <- cv.glmnet(x = dtm_train,
                      y = train$pos_neg,
                      type.measure = "auc",
                      family = "binomial",
                      alpha = 1,
                      thresh = 1e-07,
                      nfolds = 3)

coefs <- coef(model_dtm, s = "lambda.min")
coefs <- data.frame(name = coefs$dimnames[1,1][coefs$ql + 1], coefficient = coefs$sk)
```

```
print((nrow(coefs)/ncol(dtm_train)))
```

```
## [1] 0.9128412
```

```
pred_test <- predict(model_dtm, dtm_test, type= "response")[,1]
thresh <- 0.5
table(test$pos_neg, pred_test > thresh)
```

```
##
##      FALSE TRUE
## 0    1528 2145
## 1    397 15769
```

```
glmnet::auc(test$pos_neg, pred_test > thresh)
```

```
## [1] 0.6957255
```

In this one, when we added a nfolds =3 and thresh = 1e-07, we saw some improvement of 0.002 but tells us we have to increase our model iteration to better results and we should use cross validation for better AUC. Our model training model predicted the test dataset of each vocab being either a negative or positive and the best AUC we came up to is approx 70% accurate that means our model is fairly doing good

```
data <- amazon

data$text_clean <- text

skipgrams <- unnest_tokens(data, ngram, text_clean, token = "ngram", n = 9)
skipgrams$ngramID <- 1:nrow(skipgrams)
skipgrams$skipgramID <- paste(skipgrams$x_unit_id, skipgrams$ngramID, sep = "_")

head(skipgrams[, c('ngramID', 'ngram', 'skipgramID')])
```

```
##      ngramID
## 1          1
## 2          2
## 3          3
## 4          4
## 5          5
## 6          6
##
## 1          one thing book seem obvious original think however clarity
## 2          thing book seem obvious original think however clarity author
## 3          book seem obvious original think however clarity author explain
## 4          seen obvious original think however clarity author explain innovation
## 5          obvious original think however clarity author explain innovation happen
## 6          original think however clarity author explain innovation happen remarkable
##      skipgramID
## 1          -2
## 2          -3
## 3          -3
## 4          -4
## 5          -5
## 6          -6
```

```
skipgrams <- unnest_tokens(skipgrams, word, ngram)
skipgrams[skipgrams$ngramID == 1, c('skipgramID', 'word')]
```

```
##      skipgramID word
## 1          -3      one
## 2          -3      thing
## 3          -3      book
## 4          -3      seen
## 5          -3      obvious
## 6          -3      original
## 7          -3      think
## 8          -3      however
## 9          -3      clarity
```

```
library(vdify)
skipgram_probs <- pairwise_count(skipgrams, word, skipgramID, diag = T, sort = T)
```

```
## Warning: distinct() was deprecated in dplyr 0.7.0.
## Please use distinct_if() instead.
## See vignette('programming') for more help
```

```
skipgram_probs$ <- skipgram_probs$sum(skipgram_probs$)
unigram_probs <- unnest_tokens(data, word, text_clean)
unigram_probs <- count(unigram_probs, word, sort = T)
unigram_probs$ <- unigram_probs$sum(unigram_probs$)
```

```
lbound <- 30
normed_probs <- skipgram_probs[skipgram_probs$ > lbound,]
colnames(normed_probs) <- c('word1', 'word2', 'n', 'p_all', 'p_x', 'p_y')
normed_probs <- merge(normed_probs, unigram_probs[, c('word', 'p')], by.x = 'word2', by.y = 'word', all.x = T)
normed_probs <- merge(normed_probs, unigram_probs[, c('word', 'p')], by.x = 'word2', by.y = 'word', all.x = T)
```

```
head(normed_probs)
```

```
##      word1 word2 n p_all p_x p_y
## 1      18      18 105 4.727658e-07 4.286962e-06 4.286962e-06
## 2      18      cost 28 1.248815e-07 1.449912e-03 4.286962e-06
## 3      18      price 35 1.581819e-07 5.798835e-03 4.286962e-06
## 4      2      originally 25 1.119814e-07 1.381327e-04 9.400135e-04
## 5      2      idea 23 1.825813e-07 3.798659e-04 9.400135e-04
## 6      2      concern 31 1.382617e-07 2.312524e-04 9.400135e-04
```

```
# p_all = probability of seeing a given pair of words in the same window across all pairs
# p_x and p_y = probability of seeing given word across all words
normed_probs$combined <- normed_probs$p_all/normed_probs$p_x/normed_probs$p_y

normed_probs <- normed_probs[order(normed_probs$combined, decreasing = T),]
brands <- c('dryer', 'washerdryer', 'washer', 'refrigerator', 'dishwasher', 'stove')
appliances <- normed_probs[normed_probs$word1 %in% brands,]
appliances <- appliances[order(appliances$word1, appliances$combined),]
head(appliances)
```

```
##      word1 word1 word2 n p_all p_x p_y
## 236885 dishwasher 2 321 1.431878e-06 9.486135e-04 8.062328e-21
## 237357 dishwasher 3 168 7.492826e-07 5.888654e-04 8.062328e-21
## 236511 dishwasher ability 61 2.728634e-07 7.192401e-05 8.062328e-21
## 235075 dishwasher whole 569 2.537717e-06 1.141558e-03 8.062328e-21
## 237984 dishwasher absolute 25 1.115814e-07 3.548596e-05 8.062328e-21
##      p_combined
## 236885 8.6563969
## 237357 1.5458969
## 236511 1.6390182
## 235075 9.982184
## 237984 1.3549128
## 237255 0.9425169
```

```
normed_probs$mi <- log(normed_probs$combined)
pml_matrix <- cast_sparse(normed_probs, word1, word2, pml)
```

```
library(rliba)
pml_svd <- rliba(pml_matrix, 256, maxit = 1e3)
word_vectors <- pml_svd$u
rownames(word_vectors) <- rownames(pml_matrix)
```

Matching Words

```
library("tidyverse")
matching_word <- function(word_vectors, selected_vector) {
  similarities <- word_vectors %>%
    select(selected_vector) %>%
    arrange(desc(similar_prob))

  similarities %>%
    mutate(word_similar = rownames(similarities)) %>%
    select(word_similar, similar_prob)
}
```

```
dryer <- matching_word(word_vectors, word_vectors['dryer',])
rownames(dryer) <- NULL
```

```
washerdryer <- matching_word(word_vectors, word_vectors['washerdryer',])
rownames(washerdryer) <- NULL
```

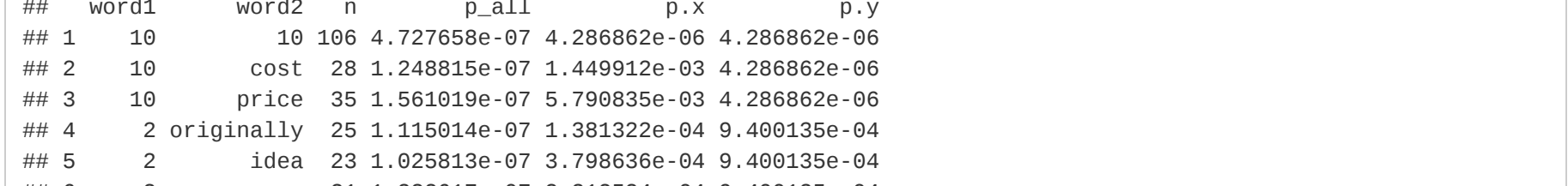
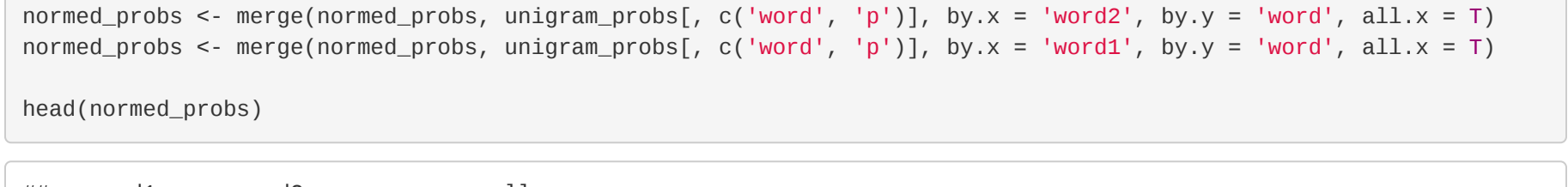
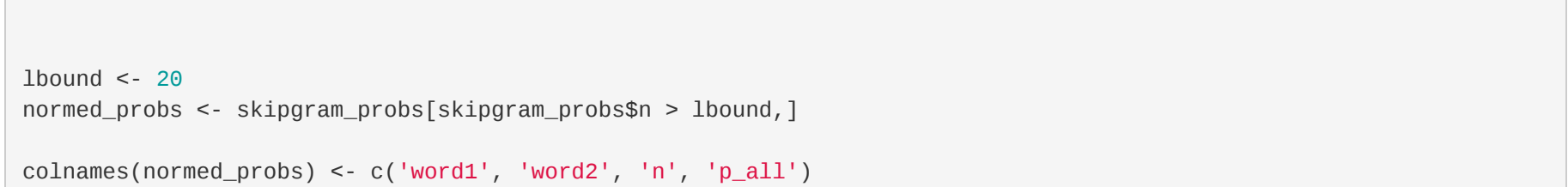
```
washer <- matching_word(word_vectors, word_vectors['washer',])
rownames(washer) <- NULL
```

```
refrigerator <- matching_word(word_vectors, word_vectors['refrigerator',])
rownames(refrigerator) <- NULL
```

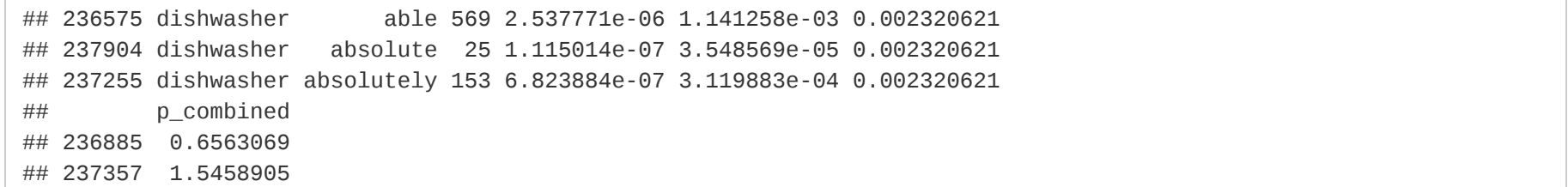
```
dishwasher <- matching_word(word_vectors, word_vectors['dishwasher',])
rownames(dishwasher) <- NULL
```

```
stove <- matching_word(word_vectors, word_vectors['stove',])
rownames(stove) <- NULL
```

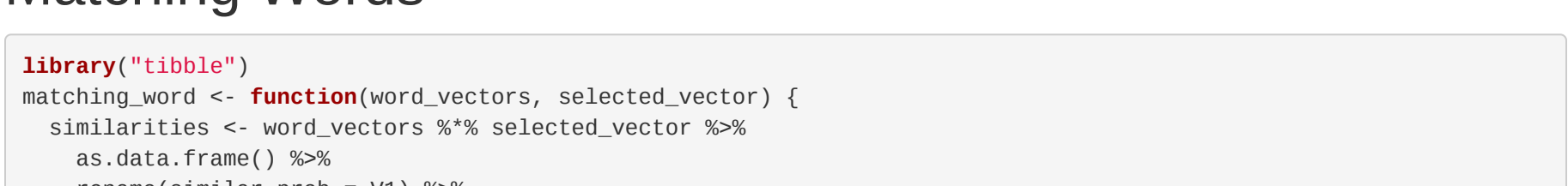
```
plot_similar <- function(df, title){
  string_title <- gsub(" ", "<div> ", title)
  ggplot(df[,2:11], aes(word_similar, similar_prob)) + geom_bar(stat='identity') + coord_flip() + theme_classic() +
    scale_y_continuous(expand = c(0,0)) +
    labs(x = NULL, title = paste1('Top 10 words from Word Vector that is associated with ', "'", string_title, "'"),
         subtitle = "Based on the Amazon Reviews, calculated using counts and matrix factorization")
}
```



We can see the word "Machine" is the most associated with washer than "dryer" and "washing".



In this graph we see how refrigerator is associated with ice as intuitive thinking and samsung as brand, and we see model as could be the model of the brand.



We can see the "range" is very high probability showing with stove and oven cooking second since they mean the s use in natural language



The advantage of the word embedding is giving you a intuitive relationships between words. We going with all the top 10 similar word are making sense. I noticed that there are some of words that can be off the relationship. Also, usually the results would give you a synonym of the same word rather than the words that come with it. I would say that word embedding can a very high performance when factor log and processing large data which needs more instances and we that can handle the computation