Group 304: R outputs Fake Job Posting Prediction

• Preprocessing of Data

i. Elimination of unused variable

```
> job_posting=job_posting[,-c(1,2,3,4,16,17)]
> colnames(job_posting)
[1] "salary_range" "company_profile" "description" "requirements" "benefits"
[6] "telecommuting" "has_company_logo" "has_questions" "employment_type" "required_experience"
[11] "required_education" "fraudulent"
> |
```

ii. Replacing empty values with 'NA'

```
> job_posting = read.csv("fake_job_postings.csv",header=T)
> job_posting[job_posting==""] <- NA
> I
```

iii. Converting nominal variable to binary variable and creating N-1 dummy variables

```
> library(dummies)
> job_posting=dummy.data.frame(job_posting,names=c("employment_type"))
Warning message:
In model.matrix.default(\sim x - 1, model.frame(\sim x - 1), contrasts = FALSE) :
 non-list contrasts argument ignored
> colnames(job_posting)
[1] "salary_range"
                                "company_profile"
                                                            "description"
                                                                                       "requirements"
[5] "benefits"
                                "telecommuting"
                                                            "has_company_logo"
                                                                                       "has_questions"
[9] "employment_typeContract" "employment_typeFull-time" "employment_typeOther"
                                                                                       "employment_typePart-time"
[13] "employment_typeTemporary" "employment_typeNA"
                                                           "required_experience"
                                                                                       "required_education"
[17] "fraudulent"
> #To drop one variable so that we have N-1 dummy variables
> job_posting=job_posting[,-c(14)]
```

iv. Assigning variables

```
> x1 <- job_posting$description
> x2 <- job_posting$telecommuting
> x3 <- job_posting$has_company_logo
> x4 <-job_posting$has_questions
> x5 <- job_posting$employment_typeContract
> x6 <- job_posting$employment_typeFull-time`
> x7 <- job_posting$employment_typeOther
> x8 <- job_posting$employment_typePart-time`
> x9 <- job_posting$employment_typeTemporary
> x10 <- job_posting$fraudulent
> |
```

v. R reads labels as nominal variables. To convert it into numeric, we use the factor function.

```
> x10f <- factor(x10)
> |
```

Two Sample Hypothesis Testing

- 1. There are 8472 job postings for which 'has_questions'=0 and are genuine.
- 2. There are 8542 job postings for which 'has_questions'=1 and are genuine.
- 3. There are 616 job postings for which 'has questions'=0 and are fake.
- 4. There are 250 job postings for which 'has_questions'=1 and are fake.

Null Hypothesis H_0 : The proportion of job postings with 'has_questions'=1 and are fake = 0.048 (Have used proportion instead of mean because it is binary data)

Alternate Hypothesis H_a: The proportion of job postings with 'has_questions'=1 and are fake is not equal to 0.048.

p-value < 0.05. Hence, null hypothesis rejected.

• KNN Model

Pre-processing of column 'description'

```
> library(tm)
> corpus_object <- Corpus(VectorSource(x1))
> corpus_object <- tm_map(corpus_object, removePunctuation)
Warning message:
In tm_map.SimpleCorpus(corpus_object, removePunctuation) :
 transformation drops documents
> corpus_object <- tm_map(corpus_object, removewords, stopwords(kind = "en"))
Warning message:
In tm_map.SimpleCorpus(corpus_object, removeWords, stopwords(kind = "en")) :
 transformation drops documents
> corpus_object <- tm_map(corpus_object, stemDocument)
Warning message:
In tm_map.SimpleCorpus(corpus_object, stemDocument) :
 transformation drops documents
> #Converting all the remaining high frequency words to a DocumentTermMatrix
> frequencies <- DocumentTermMatrix(corpus_object)
> sparse_data_desc <- removeSparseTerms(frequencies, 0.995)</pre>
> sparse_data_jpdesc <- as.data.frame(as.matrix(sparse_data_desc))</pre>
> colnames(sparse_data_jpdesc) <- make.names(colnames(sparse_data_jpdesc))
> sparse_data_jpdesc$fraudulent <- job_posting$fraudulent
> colnames(sparse_data_jpdesc) <- make.unique(colnames(sparse_data_jpdesc), sep = "_")</pre>
> set.seed(2000)
> sparse_data_jpdesc$fraudulent=factor(sparse_data_jpdesc$fraudulent)
```

For k=151

Accuracy = 95.16%

For k=71

Accuracy = 95.16%

For k=27

```
> knn.27 <- knn(train.jp,test.jp,train.fr,k=27)
> library(caret)
> confusionMatrix((table(knn.27 ,test.fr)))
Confusion Matrix and Statistics
test.fr
knn.27
                  1
     0 13610
                688
                Accuracy: 0.9518
                  95% ci : (0.9481, 0.9552)
    No Information Rate: 0.9516
    P-Value [Acc > NIR] : 0.479
                   Карра: 0.0106
 Mcnemar's Test P-Value : <2e-16
             Sensitivity: 0.99985
             Specificity: 0.00578
          Pos Pred Value : 0.95188
         Neg Pred Value : 0.66667
             Prevalence: 0.95162
   Detection Rate : 0.95148
Detection Prevalence : 0.99958
      Balanced Accuracy: 0.50282
        'Positive' Class: 0
```

Accuracy= 95.18%

For k=5

Accuracy =96.41%

• Logistic Regression Model

```
> #Logistic Regression Model
> fit <- glm(x10f ~ x2+x3+x4+x5+x6+x7+x8+x9 , data=train.data, family=binomial())
> summary(fit)
glm(formula = x10f \sim x2 + x3 + x4 + x5 + x6 + x7 + x8 + x9, family = binomial(),
       data = train.data)
Deviance Residuals:
Min 1Q Median 3Q Max
-1.0677 -0.2547 -0.2058 -0.1663 3.0339
Coefficients:
                    S:
Estimate Std. Error z value Pr(>|z|)
-1.30937 0.07883 -16.610 < 2e-16 ***
0.56391 0.14490 3.892 9.95e-05 ***
-2.10267 0.07779 -27.029 < 2e-16 ***
(Intercept) -1.30937
x2 0.56391
                                                                    < 2e-16 ***
х3
                                        0.08198 -5.251 1.51e-07 ***
0.17236 -4.349 1.37e-05 ***
0.08532 -5.065 4.09e-07 ***
0.28953 0.833 0.40484
0.14813 3.253 0.00114 **
0.71908 -2.570 0.01016 *
                    -0.43047
-0.74964
x5
x6
x7
                    -0.43213
                      0.24118
x8
                      0.48182
                    -1.84817
Signif. codes: 0 '***' 0.001 '**' 0.01 '*' 0.05 '.' 0.1 ' '1
(Dispersion parameter for binomial family taken to be 1)
Null deviance: 6933.1 on 17879 degrees of freedom Residual deviance: 5875.7 on 17871 degrees of freedom
AIC: 5893.7
Number of Fisher Scoring iterations: 7
              > prob=predict(fit,type="response",newdata=test.data)
              > prob=predict(fit,type="response",newdata=test.data)
warning message:
'newdata' had 3576 rows but variables found have 17880 rows
> #cut-off value to calculate accuracy is 0.5
> for(i in 1:length(prob)){
+    if (prob[i] > 0.5){
+       prob[i]=1}
+    else {
                   prob[i]=0
}
                 }
library(Metrics)
              > accuracy(test.label,prob)
[1] 0.9549776
```

Forward Selection Model:

```
> base=glm(x10f~x7, data=train.data, family=binomial())
> model1= step(base, scope=list(upper=fit,lower=~1),direction="both", trace=F)
> summary(model1)
glm(formula = x10f \sim x3 + x8 + x4 + x2 + x6 + x5 + x9, family = binomial(),
    data = train.data)
Deviance Residuals:
Min 1Q Median 3Q Max
-1.0673 -0.2564 -0.2059 -0.1664 3.0333
Coefficients:
             (Intercept) -1.29745
x3
x8
x4
x2
x6
x5
х9
Signif. codes: 0 '*** 0.001 '** 0.01 '* 0.05 '.' 0.1 ' '1
(Dispersion parameter for binomial family taken to be 1)
Null deviance: 6933.1 on 17879 degrees of freedom
Residual deviance: 5876.3 on 17872 degrees of freedom
AIC: 5892.3
Number of Fisher Scoring iterations: 7
```

```
> prob=predict(model1,type="response",newdata=test.data)
warning message:
'newdata' had 3576 rows but variables found have 17880 rows
> #cut-off value to calculate accuracy is 0.5
> for(i in 1:length(prob)){
+    if (prob[i] > 0.5){
+       prob[i]=1}
+    else {
+       prob[i]=0
+    }
+ }
> library(Metrics)
> accuracy(test.label,prob)
[1] 0.9549776
```

Backward Selection Model:

```
> model2= step(fit,direction="backward",trace=F)
> summary(model2)
call:
glm(formula = x10f \sim x2 + x3 + x4 + x5 + x6 + x8 + x9, family = binomial()
   data = train.data)
Deviance Residuals:
Min 1Q Median 3Q Max
-1.0673 -0.2564 -0.2059 -0.1664 3.0333
Coefficients:
Signif. codes: 0 '***' 0.001 '**' 0.01 '*' 0.05 '.' 0.1 ' ' 1
(Dispersion parameter for binomial family taken to be 1)
Null deviance: 6933.1 on 17879 degrees of freedom
Residual deviance: 5876.3 on 17872 degrees of freedom
AIC: 5892.3
Number of Fisher Scoring iterations: 7
      > prob=predict(model2,type="response",newdata=test.data)
      prob[i]=1}
else {
         prob[i]=0
      + }
+ }
> library(Metrics)
      > accuracy(test.label,prob)
[1] 0.9549776
```

SVM Model

Accuracy = 85.156%

Feature Extraction

Code execution on the genuine post description:

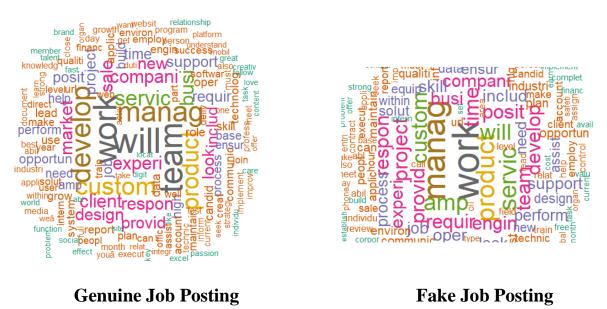
```
package 'dplyr' was built under R version 3.6.3
> description <- job_posting$description
> description <- as.character(description)</pre>
  > Genuine_description <- filter(job_posting,fraudulent==0)
 > description_G <- as.data.frame(Genuine_description$description)
> count <- seq(1, nrow(description_G), 1)
> lowerC_description_G <- sapply(count, function(c){tolower(description_G[c, 1])})</pre>
  > library(stringr)
 Warning message:
package 'stringr' was built under R version 3.6.3
  belong in the strict of the strict of the strict of the striction beautiful or striction of the striction of
  > library(tm)
 Loading required package: NLP
  Warning message:
package 'tm' was built under R version 3.6.3
> del_special_chars_description_G <- VCorpus(VectorSource(del_special_chars_description_G))
> del_special_chars_description_G <- tm_map(del_special_chars_description_G, removeWords, stopwords(kind = "en"))
> del_special_chars_description_G <- tm_map(del_special_chars_description_G, removePunctuation)</pre>
  > white_space_cleanup_description_G <- tm_map(del_special_chars_description_G, stripWhitespace)</pre>
  > library(Snowballc)
"addit"
                                                                                                                                                                                                                                                                                     "across"
                                                                                                                                                                                                                                                                                  "amp"
"area"
                                                                                                                                                                                                                                                                                                                                      "analysi"
                                                                                                                                                                                                                                 "also"
                                                                                                                                                                                                                                                                                                                                                                                           "analyt"
                                                                                                                                                                                                                               "appropri"
"back"
                                                                                                                                                                                                                                                                                                                                    "around"
"benefit"
                                                                                                                                                                                                                                                                                                                                                                                             "assign'
                                                                                                                                                                                                                                                                                    "base"
"call"
                                                                                                                                                                                                                                                                                                                                                                                           "best
                                                                                                                                                                                                                                "busi"
                                                                                                                                                                                                                                                                                                                                      "campaign"
```

Creation of word cloud:

Code execution on the fake job post description:

```
Fake_description <- filter(job_posting,fraudulent==1)
description_F <- as.data.frame(Fake_descriptionSdescription)
count <- seq(1, nrow(description_F), 1)
lowerC_description_F <- sapply(count, function(c){tolower(description_F[c, 1])})
library(stringr)
del_special_chars_description_F <- sapply(count, function(1){str_replace_all(lowerC_description_F[1], "[^[:alnum:]]", " ")})
library(stringr)</pre>
       del_special_chars_description_F <- sapply(count, function(!){str_replace_all(lowerC_description_F[!], "[A[:alnurlibrary(tm)]
del_special_chars_description_F <- VCorpus(VectorSource(del_special_chars_description_F))
del_special_chars_description_F <- tm_map(del_special_chars_description_F, removewords, stopwords(kind = "en"))
del_special_chars_description_F <- tm_map(del_special_chars_description_F, removewords)
white_space_cleanup_description_F <- tm_map(del_special_chars_description_F, stripwhitespace)
stemming_description_F <- tm_map(white_space_cleanup_description_F, stembocument)
text_description_F <- tm_map(stemming_description_F, stembocument, language = "english")
library(snowballc)
description_F <- TermDocumentMatrix(text_description_F)</pre>
> library(snowballC)
> docterm_corpus_description_F <- TermDocumentMatrix(text_description_F)
> review_Fake = removeSparseTerms(docterm_corpus_description_F, 0.99)
> findFreqTerms(review_Fake, 1000)
[1] "manag" "work"
> description_matrix_F <- as.matrix(review_Fake)
> v_fake <- sort(rowSums(description_matrix_F),decreasing=TRUE)
> d_fake <- data.frame(word = names(v_fake),freq=v_fake)
> set_seed(2020)
```

Below is the snapshot of the word cloud output for genuine and fake job postings.



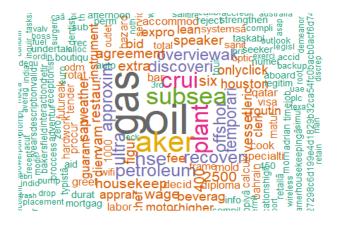
Genuine Job Posting

Fake Job Posting

Our next step would be to find the words that are present only in the fraudulent job postings and not in the genuine job postings.

Below is the snapshot of the code used for the same purpose.

The word cloud of the words that are present in fake job postings but not in the genuine job postings is given below:



From the above output, we can say that the words gas, oil, crui, aker, subsea, plant, offshor, temporari are present in the description of the fake job postings but not in genuine job postings.