

# DSP\_word2vec\_ds2

Melanie Weissenboeck

2022-11-28

## Laden von Bibliotheken und Daten

```
library("xlsx")
library(word2vec)
library(udpipe)
library(SnowballC)
library(ggplot2)
library(tm)
library(wordcloud)
library(tidytext)
library(tidyr)
library(mlbench)
library(e1071)
library(caret)
library(class)
```

## Vorverarbeiten der Texte

### Bereinigen der Texte

```
# Funktion fuer Text Bereinigung
ds2$ANF_BESCHREIBUNG <- txt_clean_word2vec(ds2$ANF_BESCHREIBUNG,
                                           ascii = FALSE, alpha = TRUE,
                                           tolower = TRUE, trim = TRUE)
```

### Erstellen einer Worteinbettung

```
# Modell trainieren fuer Einbettung
model_ds2 <- word2vec(ds2$ANF_BESCHREIBUNG, dim = 10, iter = 15)
embedding_ds2 <- as.matrix(model_ds2)
```

```
# Dimension der Einbettung
dim(embedding_ds2)
```

```
## [1] 5305 10
```

### Generieren von numerischen Prädiktoren

```
# aufteilen der Texte in einzelne Token
ds2$token <- tokenizers::tokenize_words(ds2$ANF_BESCHREIBUNG)
```

```
# Vektor der Laenge 10 fuer jedes Dokument
features2 <- matrix(nrow = 0, ncol = 10)
for (i in (1:length(ds2$ANF_BESCHREIBUNG))){
  vec_doc1 <- doc2vec(model_ds2, ds2$token[1][[1]][i], split = " ")
  features2 <- rbind(features2, vec_doc1)
}
```

## Zusammenführen mit anderen Prädiktoren

```
features2 <- as.data.frame(features2)
ds2_all <- cbind(ds2, features2)
ds2_all <- as.data.frame(ds2_all)
```

```
df <- ds2_all[, c(6, 7, 8, 9, 11, 12, 13, 14, 16, 17, 18, 19, 20)]
df[is.na(df)] <- 0
head(df)
```

```
##   ANF_RISIKO TF_ABDECKUNG AKT_RES_STATUS AKT_RES_RELEASE      V1      V2
## 1   mittel      50.0          OK          22.30 -0.6890023 -0.3156689
## 2   mittel      50.0          OK          22.30  0.0000000  0.0000000
## 3   mittel      50.0          OK          22.20  0.0000000  0.0000000
## 4   mittel      50.0          OK          22.20  0.7015501 -0.3013070
## 5   gering       0.0          OK           21x -2.1290462  0.1590461
## 6   mittel       0.8        FAILED          22.30 -1.1517696  0.9719147
##           V3           V4           V6           V7           V8           V9           V10
## 1  1.1615236  0.54112215 -1.71687869  0.2562611 -2.0580583 -0.2272617 -0.5302653
## 2  0.0000000  0.00000000  0.00000000  0.0000000  0.0000000  0.0000000  0.0000000
## 3  0.0000000  0.00000000  0.00000000  0.0000000  0.0000000  0.0000000  0.0000000
## 4 -0.7693904 -0.02538728  1.56988087  1.8889872 -0.7688528  1.3791198 -0.4377108
## 5 -1.2191108 -0.03083323 -0.05161381  0.1990778 -0.9227785  0.7190004 -1.3394860
## 6  0.5045080 -0.43781123  0.53666181  2.0884791  0.1566678 -1.3472567 -0.8906209
```

## Normalisieren numerischer Spalten

```
set.seed(1234)

# definiere normalisierungsfunktion
min_max_norm <- function(x) {
  (x - min(x)) / (max(x) - min(x))
}

# alle spalten normalisieren
df[, 5:13] <- as.data.frame(lapply(df[, 5:13], min_max_norm))
df[2] <- as.data.frame(lapply(df[2], min_max_norm))

df$ANF_RISIKO <- as.factor(df$ANF_RISIKO)
df$AKT_RES_STATUS <- as.factor(df$AKT_RES_STATUS)
df$AKT_RES_RELEASE <- as.factor(df$AKT_RES_RELEASE)
summary(df)
```

```
##   ANF_RISIKO   TF_ABDECKUNG   AKT_RES_STATUS AKT_RES_RELEASE
## gering: 633   Min.      :0.00000   FAILED: 577   21x      :1146
## hoch :1122   1st Qu.:0.04121   OK      :2363   22.10    : 434
## mittel:1366   Median :0.17229   OPEN   : 181   22.20    : 727
##           Mean      :0.30279           22.30    : 326
```

```
##          3rd Qu.:0.50348          OLDERT21: 488
##          Max.    :1.00000
##          V1          V2          V3          V4
## Min.    :0.0000    Min.    :0.0000    Min.    :0.0000    Min.    :0.0000
## 1st Qu.:0.5137    1st Qu.:0.4822    1st Qu.:0.3991    1st Qu.:0.4306
## Median :0.5137    Median :0.4822    Median :0.3991    Median :0.4306
## Mean    :0.5150    Mean    :0.4869    Mean    :0.4011    Mean    :0.4343
## 3rd Qu.:0.5137    3rd Qu.:0.4822    3rd Qu.:0.3991    3rd Qu.:0.4306
## Max.    :1.0000    Max.    :1.0000    Max.    :1.0000    Max.    :1.0000
##          V6          V7          V8          V9
## Min.    :0.0000    Min.    :0.0000    Min.    :0.0000    Min.    :0.0000
## 1st Qu.:0.4708    1st Qu.:0.3885    1st Qu.:0.4827    1st Qu.:0.6408
## Median :0.4708    Median :0.3885    Median :0.4827    Median :0.6408
## Mean    :0.4723    Mean    :0.3890    Mean    :0.4844    Mean    :0.6361
## 3rd Qu.:0.4708    3rd Qu.:0.3885    3rd Qu.:0.4827    3rd Qu.:0.6408
## Max.    :1.0000    Max.    :1.0000    Max.    :1.0000    Max.    :1.0000
##          V10
## Min.    :0.0000
## 1st Qu.:0.5736
## Median :0.5736
## Mean    :0.5697
## 3rd Qu.:0.5736
## Max.    :1.0000
```

## Klassifikation

### Erstellen von Train- / Test-Split

```
# partition erstellen
part <- createDataPartition(df$ANF_RISIKO, times = 1, p = 0.80)
# extract training set
X_train <- df[part$Resample1, ]
# extract testing set
X_test <- df[-part$Resample1, ]
# extract target
y_train <- df[part$Resample1, 1]
y_test <- df[-part$Resample1, 1]
```

### NaiveBayes Klassifikation

```
model_nb = naiveBayes(ANF_RISIKO ~ ., data = X_train)
```

```
pred_nb <- predict(model_nb, X_test)
mat.nb <- confusionMatrix(pred_nb, X_test$ANF_RISIKO, mode = "prec_recall")
mat.nb
```

```
## Confusion Matrix and Statistics
##
##          Reference
## Prediction gering hoch mittel
##      gering      13   11     3
##      hoch       73  193    231
##      mittel     40   20    39
```

```
##
## Overall Statistics
##
##           Accuracy : 0.3933
##           95% CI : (0.3547, 0.4329)
##       No Information Rate : 0.4382
##       P-Value [Acc > NIR] : 0.9896
##
##           Kappa : 0.0442
##
## Mcnemar's Test P-Value : <2e-16
##
## Statistics by Class:
##
##           Class: gering Class: hoch Class: mittel
## Precision           0.48148      0.3883      0.3939
## Recall              0.10317      0.8616      0.1429
## F1                  0.16993      0.5354      0.2097
## Prevalence          0.20225      0.3596      0.4382
## Detection Rate      0.02087      0.3098      0.0626
## Detection Prevalence 0.04334      0.7978      0.1589
## Balanced Accuracy    0.53750      0.5499      0.4857
```

## KNN Klassifikation

```
model_knn <- train(ANF_RISIKO ~ ., data = X_train, "knn",
                  trControl = trainControl(method = "cv", number = 5))

pred_knn <- predict(model_knn, X_test, type = "raw")
mat.knn <- confusionMatrix(pred_knn, X_test$ANF_RISIKO, mode = "prec_recall")
mat.knn
```

```
## Confusion Matrix and Statistics
##
##           Reference
## Prediction gering hoch mittel
##      gering      46   15    27
##      hoch       10  120    35
##      mittel     70   89   211
##
## Overall Statistics
##
##           Accuracy : 0.6051
##           95% CI : (0.5655, 0.6437)
##       No Information Rate : 0.4382
##       P-Value [Acc > NIR] : < 2.2e-16
##
##           Kappa : 0.3589
##
## Mcnemar's Test P-Value : 1.855e-09
##
## Statistics by Class:
##
##           Class: gering Class: hoch Class: mittel
```

## Precision	0.52273	0.7273	0.5703
## Recall	0.36508	0.5357	0.7729
## F1	0.42991	0.6170	0.6563
## Prevalence	0.20225	0.3596	0.4382
## Detection Rate	0.07384	0.1926	0.3387
## Detection Prevalence	0.14125	0.2648	0.5939
## Balanced Accuracy	0.64029	0.7115	0.6593