

Opis problemu

- Zbiór danych: Real or Not? NLP with Disaster Tweets (kaggle.com)
- Dwa zbiory: treningowy (7163 rekordów) i testowy (3263 rekordy)
- Dane: pliki csv [id, location, keyword, text, (target)]
- Zmienna celu: tweet opisuje faktyczne zdarzenie katastroficzne (tak bądź nie)
- Opis problemu:
 Zbudowanie modelu na podstawie danych treningowych, który
 będzie klasyfikował czy dany tweet opisuje faktyczną katastrofę czy nie





Przyjrzyjmy się danym

Rzut okiem na pierwsze wiersze:



- Przemyślenia:
 - Czy id będzie nam w jakikolwiek sposób potrzebne?
 - Mamy kolumny, w których brakuje danych co z nimi?
- Akcje:
 - Usuwamy Id
 - Sprawdzamy co zrobić z [location] i [keyword]

```
> missing_train_data <- colsums(sapply(train_data, is.na))
> missing_train_data
    id keyword location text target
        0 61 2533 0 0
```

- Dużo braków w location
- Co zrobić z keyword?

Przyjrzyjmy się danym

Zobaczmy jak procentowo do całego zbioru mają się nasze braki

- Ponad 30% braków w kolumnie location -> Usuwam kolumnę
- Ocena przydatności kolumny keyword względem zmiennej celu -> wydają się istotne, postanawiam je pozostawić do dalszej analizy

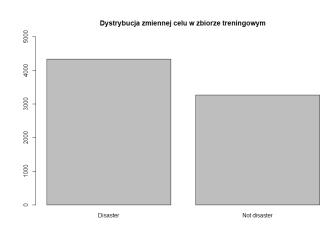
```
> library(dplyr)
> train_data$keyword <- sapply(train_data$keyword, as.character)
> train_data$keyword[is.na(train_data$keyword)] <- "no category"
>
> by_target <- train_data %>% group_by(keyword) %>% summarize(target_mean = sum(target)/n()) %>% arrange(desc(target_mean))
> View(by_target)
```

keyword ‡	target_mean ‡
debris	1.00000000
derailment	1.00000000
wreckage	1.00000000
outbreak	0.97500000
oil%20spill	0.97368421
typhoon	0.97368421
suicide%20bombing	0.96969697
suicide%20bomber	0.96774194
bombing	0.93103448
rescuers	0.91428571
suicide%20bomb	0.91428571
nuclear%20disaster	0.91176471
evacuated	0.88888889

	keyword ‡	target_mean ^
	aftershock	0.00000000
	body%20bags	0.02439024
	ruin	0.02702703
	blazing	0.02941176
	body%20bag	0.03030303
	electrocute	0.03125000
	screaming	0.05555556
	traumatised	0.05714286
	blew%20up	0.06060606
10	panicking	0.06060606
	blight	0.06250000
	wrecked	0.07692308
	explode	0.07894737

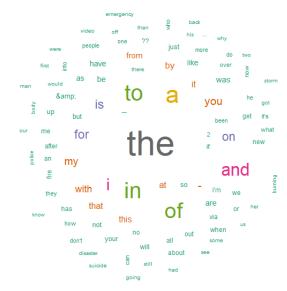
Wizualizacje

Dystrybucja zmiennej celu w zbiorze treningowym



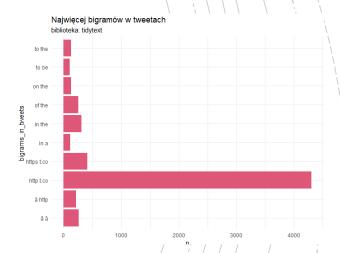
 Dystrybucja wydaje się być zrównoważona, nie widzę konieczności redystrybucji danych do nauki

Chmura słów – treść tweet-a



 Najczęstsze słowa to angielskie stop-words – koniecznie trzeba oczyścić dane

Bigramy – analiza występujących par słów



 Dzięki bigramom widzimy pary słów.
 Wskazują, że do oczyszczenia są również chociażby linki

Połączenie zbiorów w jeden

Wszystkie operacje
 procesowania tekstu takie
 same dla zbioru treningowego
 i testowego

```
all_tweets_df <- bind_rows(train_data, test_data)
all_data_corpus <- Corpus(VectorSource(all_tweets_df$text))</pre>
```



2

Lowercase dla treści tweetów

```
all_data_corpus <- tm_map(
   all_data_corpus, content_transformer(
   stri_trans_tolower))</pre>
```

4

Usunięcie nazw użytkowników

```
# Usuniecie nazw użyttkowników
remove_usernames <- function(x){
   gsub("@[^[:space:]]*", "", x)
}
all_data_corpus <- tm_map(
   all_data_corpus, content_transformer(
   remove_usernames))</pre>
```

3

Usunięcie url-ów z tweetów

```
# Ususniecie linków (ulr-ów)
remove_urls <- function(x) {
   gsub("http[^[:space:]]*", "", x)
}
all_data_corpus <- tm_map(
   all_data_corpus, content_transformer(
   remove_urls))</pre>
```

5

Usunięcie słów jedno i dwu literowych

```
# Usuniecie słów jedno i dwu wyrazowych
remove_one_char <- function(x){
    gsub(" . ", " ", x)
}
remove_two_chars <- function(x){
    gsub(" . . ", " ", x)
}
all_data_corpus <- tm_map(all_data_corpus, content_transformer(remove_one_char))
all_data_corpus <- tm_map(all_data_corpus, content_transformer(remove_two_chars))
# stemming słów z textów tweetów
all_data_corpus <- tm_map(
    all_data_corpus, stemDocument, "english")</pre>
```

Usunięcie znaków innych niż litery

```
remove_anything_but_text <- function(x){
   gsub("[^[:alpha:][:space:]]*", "", x)
}
all_data_corpus <- tm_map(
   all_data_corpus, content_transformer(
    remove_anything_but_text))</pre>
```



Usunięcie angielskich stop-words

```
#install.packages('qdap')
library(qdap)
# Usuniecie angielskich stop-words
all_data_corpus <- tm_map(
   all_data_corpus, removewords, Top200words)</pre>
```



Stemming

@Username

```
# stemming słów z textów tweetów
all_data_corpus <- tm_map(
   all_data_corpus, stemDocument, "english")</pre>
```

Przekształcony tekst tweet-a



Name

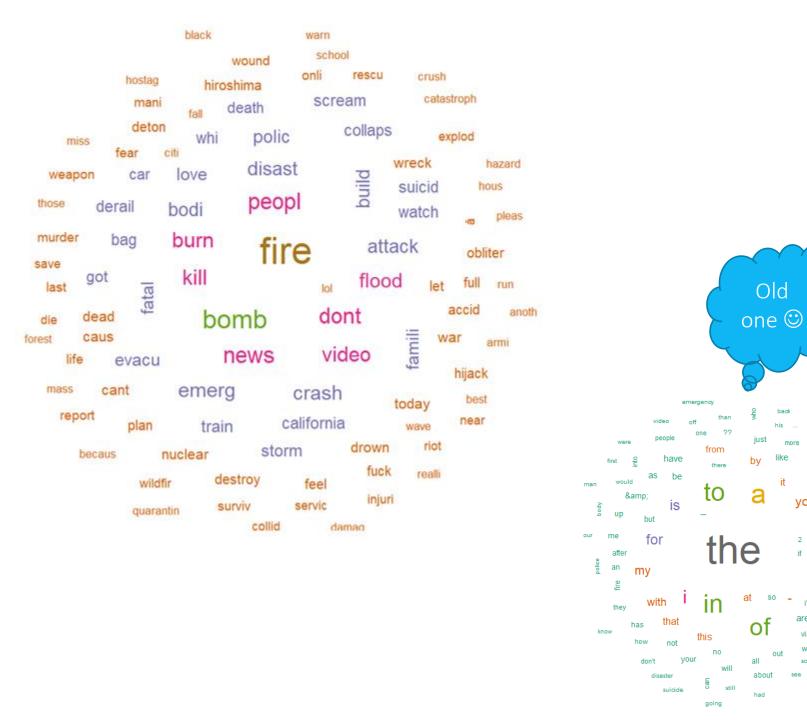
Follow

'CALIFORNIA IS BURNING:' Gov. Jerry Brown told reporters at a press conference that California is experiencing... http://t.co/arzeMSR7FQ



california burn gov jerri brown told report press confer california experienc

Procesowanie tekstu





you

on

and

Budowanie modelu

Budowa DocumentTermMatrix

DTM <- DocumentTermMatrix(train_data_corpus)
inspect(DTM[1:10, 50:60])</pre>

```
Non-/sparse entries: 12/98

Sparsity : 89%

Maximal term length: 7

Weighting : term frequency (tf)
```

		_								
Docs	across	afraid	build	emerg	happen	hill	spring	street	top	wood
1	0	0	0	Ō	0	0	Ō	0	0	0
10	0	1	0	0	0	0	0	0	0	0
2	0	0	0	0	0	0	0	0	0	0
3	0	0	0	0	0	0	0	0	0	0
4	0	0	0	0	0	0	0	0	0	0
5	0	0	0	0	0	0	0	0	0	0
6	0	0	0	0	0	0	0	0	0	0
7	0	0	0	0	0	0	1	1	0	0
8	0	0	0	0	0	1	0	0	1	1
9	1	0	1	1	1	0	0	1	0	0

Usunięcie rzadkich słów

```
sparse_DTM <- removeSparseTerms(DTM, 0.995)
tweetsSparse <- as.data.frame(as.matrix(sparse_DTM))
colnames(tweetsSparse) <- make.names(colnames(tweetsSparse))
tweetsSparse$target <- train_data$target</pre>
```

Podział na zbiór treningowy i testowy

```
# nauka - podział na 70% danych treningowych i 30 testowych
smp_size <- floor(0.7 * nrow(tweetsSparse))
set.seed(123)
train_ind <- sample(seq_len(nrow(tweetsSparse)), size = smp_size)
train <- tweetsSparse[train_ind, ]
test <- tweetsSparse[-train_ind, ]</pre>
```

Drzewo decyzyjne

Trenowanie modelu + ewaluacja

```
# ternowanie modelu - drzewo decyzyjne
library('rpart')

tweetCART <- rpart(target ~ . , data = train, method = "class")

predictCART <- predict(tweetCART, newdata = test, type = "class")

cmat_CART <- table(test$target, predictCART)

cmat_CART
accu_CART <- (cmat_CART[1,1] + cmat_CART[2,2])/sum(cmat_CART)
accu_CART</pre>
```



SVM

Trenowanie modelu + ewaluacja





Co zostało zrobione

- Analiza zbioru danych
- Procesowanie tekstu celem przygotowania do analizy
- Wizualizacja i analiza uzyskanych rezultatów po procesowaniu
- Trenowanie modelu i ewaluacja

Dalsze plany

- Próba poprawienia wyniku przewidywania przez model użycie bardziej zaawansowanych algorytmów, np. xgboost
- Praca nad doborem hiperparametrów do trenowania analiza porównawcza wyników
- Dalsza analiza z wykorzystaniem zbioru testowego celem ukończenia konkursu na stronie Kaggle