

Supervised learning aggregated

Cornelius Erfort

5/10/2021

Contents

1	Setting up	1
1.1	Loading packages	1
1.2	Loading data	2
1.3	Merging categories	2
1.4	Creating the document frequency matrix (dfm)	3
2	Textmodels	3
2.1	Multinomial Naive Bayes classification model	3
2.2	SVM	4
2.3	Next algorithm...	5
3	Ensemble methods: SuperLearner	5
4	Classification of unlabelled data	7
4.1	Using the NB textmodel	7
4.2	Aggregation of the issues categories over time and party	8
4.3	Plotting issue attention: Examples	9

1 Setting up

This script requires the files “sample_germany.dta” and “data_joint.RDS” in the parent directory. It also writes and reads two csv-files from the parent directory.

1.1 Loading packages

This script is based mainly on the functions of the quanteda package. For the cross-validation of the textmodels, quanteda.classifiers has to be loaded from GitHub.

```
packages <- c("quanteda", "quanteda.textmodels", "dplyr", "caret", "randomForest",
             "tm", "beepR", "rmarkdown", "e1071", "penalized", "plyr", "readr", "repr", "ggplot2",
             "rsample", "remotes", "stringr", "formatR", "readstata13", "lubridate", "reticulate")

# Using readstata13 because haven causes crash

lapply(packages[!(packages %in% rownames(installed.packages()))], install.packages)

if (!("quanteda.classifiers" %in% rownames(installed.packages()))) {
  remotes::install_github("quanteda/quanteda.classifiers")
}
```

```
invisible(lapply(c(packages, "quantda.classifiers"), require, character.only = T))
```

1.2 Loading data

The sample data for Germany consists of 2,742 labelled press releases. The dataset is not on GitHub and is loaded from the parent directory here.

```
sample_germany <- read.dta13("../sample_germany.dta", convert.factors = F)
```

```
# Correcting classification for three documents
sample_germany$issue[sample_germany$id == 229] <- 191
sample_germany$issue[sample_germany$id == 731] <- 7
sample_germany$issue[sample_germany$id == 902] <- 10
```

```
# Subset to relevant vars
germany_textpress <- sample_germany %>%
  select("header", "text", "issue", "position", "id")
```

```
# Distribution of issues in the hand-coded sample
table(germany_textpress$issue)
```

```
##
##   1   2   3   4   5   6   7   8   9  10  12  13  14  15  16  17  18  20  23  98
## 175 181 119  99 167 137  84 105 131  74 195 104  32 168 121  68  27  97  19  91
##   99 191 192
##   46 350 152
```

1.3 Merging categories

In order to improve the classification, similar topics are merged or subsumed under the “Other” category. In practice, press releases regarding, for instance, Environment and Energy are often not distinguishable. Furthermore, small categories with very few observations are not suitable for automated classification.

```
germany_textpress$issue_r1 <- as.numeric(germany_textpress$issue)
```

```
germany_textpress <- germany_textpress %>% mutate(issue_r1 = recode(issue_r1,
  `8` = 7, # Environment & Energy
  `13` = 10, # Transportation & Welfare
  `14` = 10, # Housing & Welfare
  `18` = 15, # Foreign Trade and Domestic Commerce
  `98` = 99, # Non-thematic & Other
  `23` = 99) # Culture: Too few observations
)
```

```
# Category descriptions
```

```
issue_categories <- data.frame(issue_r1 = c(1:7, 9:10, 12, 15:17, 20, 99, 191:192), issue_r1_descr = c(
```

```
# Distribution with merged categories
table(germany_textpress$issue_r1)
```

```
##
##   1   2   3   4   5   6   7   9  10  12  15  16  17  20  99 191 192
## 175 181 119  99 167 137 189 131 210 195 195 121  68  97 156 350 152
```

1.4 Creating the document frequency matrix (dfm)

We create a text corpus based on the header and text of each press release. We draw a random sample from the corpus to create a training and a test dataset. The test dataset consists of approx. one fifth of the documents.

Subsequently, we follow standard procedures for the preparation of the document frequency matrix. First, we remove stopwords and stem the words in order to better capture the similarities across documents. Second, we remove all punctuation, numbers, symbols and URLs. In a last step, we remove all words occurring in less than 0.5% or more than 90% of documents.

```
corp_press <- str_c(germany_textpress$header, " ", germany_textpress$text) %>% corpus()
```

```
## Warning: NA is replaced by empty string
```

```
# Add id var to corpus
docvars(corp_press, "id") <- germany_textpress$id
docvars(corp_press, "issue_r1") <- germany_textpress$issue_r1

# Create random sample for test dataset (size: 1/5 of all classified documents)
set.seed(300)
id_test <- sample(docvars(corp_press, "id"),
                  round(length(docvars(corp_press, "id"))/5, 0), replace = FALSE)

# Create dfm
dfmat <- corpus_subset(corp_press) %>%
  dfm(remove = stopwords("de"), # Stem and remove stopwords, punctuation etc.
      stem = T,
      remove_punct = T,
      remove_number = T,
      remove_symbols = T,
      remove_url = T) %>%
  dfm_trim(min_docfreq = 0.005, # Remove words occurring <.5% or > 80% of docs
          max_docfreq = .9,
          docfreq_type = "prop") %>%
  suppressWarnings()

# Create training and test set
dfmat_training <- dfm_subset(dfmat, !(id %in% id_test))
dfmat_test <- dfm_subset(dfmat, id %in% id_test)

# Make order of documents random
dfmat <- dfmat[sample(1:ndoc(dfmat), ndoc(dfmat)), ]
```

2 Textmodels

2.1 Multinomial Naive Bayes classification model

We calculate a Multinomial Naive Bayes text classification model. Multinomial NB models take into account the number of times a word occurs in a document, whereas Bernoulli NB models use the presence or absence of words only.

```
# Five-fold cross-validation for every possible parameter combination
nb_eval <- textmodel_evaluate(dfmat, dfmat$issue_r1, k = 5, model = "textmodel_nb",
                             fun = c("accuracy", "precision", "recall", "f1_score"), parameters = list(prior = c("uniform",
                                                         "docfreq", "termfreq"), distribution = c("multinomial", "Bernoulli"), smooth = c(1,
```

```

    2, 3)))
head(nb_eval)

##    k accuracy precision    recall  f1_score  prior distribution smooth time
## 1 1 0.6065574 0.5992441 0.5873313 0.5863247 uniform multinomial      1 0.30
## 2 2 0.6496350 0.6444556 0.6412655 0.6361559 uniform multinomial      1 0.28
## 3 3 0.6605839 0.6420788 0.6443840 0.6370078 uniform multinomial      1 0.31
## 4 4 0.6514599 0.6253740 0.6326816 0.6226715 uniform multinomial      1 0.26
## 5 5 0.5974499 0.5760496 0.5876512 0.5706353 uniform multinomial      1 0.27
##          seed
## 1 1621245204
## 2 1621245204
## 3 1621245204
## 4 1621245204
## 5 1621245204

aggregate(cbind(accuracy, precision, recall, f1_score, time, seed) ~ prior + distribution +
  smooth, nb_eval[, -c(1)], mean) %>%
  arrange(desc(accuracy))

##      prior distribution smooth accuracy precision    recall  f1_score   time
## 1  docfreq multinomial      1 0.6338665 0.6217454 0.6165320 0.6102681 0.5820
## 2  uniform multinomial      1 0.6331372 0.6174404 0.6186627 0.6105591 0.2840
## 3  termfreq multinomial      1 0.6324093 0.6174881 0.6150507 0.6081595 0.7140
## 4  termfreq multinomial      2 0.6221910 0.6111748 0.5924094 0.5888383 0.3520
## 5  uniform multinomial      2 0.6214637 0.6119762 0.5952630 0.5917254 0.6740
## 6  docfreq multinomial      2 0.6214617 0.6130568 0.5913348 0.5887123 0.3940
## 7  uniform multinomial      3 0.6156303 0.6129651 0.5783785 0.5762503 0.2920
## 8  termfreq multinomial      3 0.6127126 0.6115502 0.5727716 0.5706607 0.6700
## 9  docfreq multinomial      3 0.6123476 0.6066068 0.5704696 0.5682803 0.3580
## 10 uniform Bernoulli      1 0.5481092 0.5893067 0.4942481 0.5004393 0.8150
## 11 termfreq Bernoulli      1 0.5453728 0.5837934 0.4885404 0.4934106 0.6175
## 12 docfreq Bernoulli      1 0.5435497 0.5830665 0.4871600 0.4922446 0.6500
##          seed
## 1 1621245204
## 2 1621245204
## 3 1621245204
## 4 1621245204
## 5 1621245204
## 6 1621245204
## 7 1621245204
## 8 1621245204
## 9 1621245204
## 10 1621245204
## 11 1621245204
## 12 1621245204

```

Assuming a multinomial distribution of text features leads to a higher accuracy of the models compared to a Bernoulli distribution. The other benchmark parameters confirm this finding.

There is no clear pattern in regard to the effect of the priors on the quality of the models. A smoothing parameter of 1 for the feature counts seems optimal.

2.2 SVM

Linear predictive models estimation based on the LIBLINEAR C/C++ Library.

(Not running here to save time.)

```
# Five-fold cross-validation for every possible parameter combination sum_eval <-  
# textmodel_evaluate(dfmat, dfmat$issue_r1, k = 5, model = 'textmodel_sum', fun =  
# c('accuracy', 'precision', 'recall', 'f1_score'), parameters = list(weight =  
# c('uniform', 'docfreq', 'termfreq'), type = c(0:7))) head(sum_eval)  
# aggregate(cbind(accuracy, precision, recall, f1_score, time, seed) ~ weight +  
# type, sum_eval[, -c(1)], mean) %>% arrange(desc(accuracy))
```

None of the configurations lead to a higher accuracy compared to the benchmark NB model.

Regarding the type of linear models, type 0 (L2-regularized logistic regression, primal) and 7 (L2-regularized logistic regression, dual) seem to provide optimal results. There is no clear pattern regarding the weights.

The calculation of the models requires much more time/computing compared to the NB models.

2.3 Next algorithm...

```
# sum_eval <- textmodel_evaluate(dfmat, dfmat$issue_r1, k = 5, model =  
# 'textmodel_sum', fun = c('accuracy', 'precision', 'recall', 'f1_score'),  
# parameters = list(weight = c('uniform', 'docfreq', 'termfreq')) head(sum_eval)  
# aggregate(cbind(accuracy, time, seed) ~ weight, sum_eval[, -c(1)], mean) %>%  
# arrange(desc(accuracy))
```

3 Ensemble methods: SuperLearner

Prepare the datasets for the python code chunk.

```
# Get vector indicating training sample  
training <- !(1:ndoc(corp_press) %in% id_test)  
  
# Subset training sample and write  
if (!file.exists("../train_data.csv")) as.data.frame(as.matrix(dfm_trim(dfmat, min_docfreq = 15,  
  max_docfreq = 0.8 * nrow(dfmat), verbose = T))) %>%  
  write.csv(train_data, "../train_data.csv")  
  
# Get labels for training sample and write  
if (!file.exists("../train_labels.csv")) write.csv(dfmat$issue_r1, "../train_labels.csv")
```

Run the SuperLearner in Python.

```
# Load packages  
import pandas as pd  
import numpy as np  
  
# Load sklearn tools  
from sklearn.model_selection import train_test_split  
from sklearn.metrics import accuracy_score  
  
# Load classifiers  
from sklearn.neighbors import KNeighborsClassifier  
from sklearn.linear_model import LogisticRegression  
from sklearn.tree import DecisionTreeClassifier  
from sklearn.svm import SVC  
from sklearn.naive_bayes import GaussianNB  
from sklearn.ensemble import AdaBoostClassifier
```

```

from sklearn.ensemble import BaggingClassifier
from sklearn.ensemble import RandomForestClassifier
from sklearn.ensemble import ExtraTreesClassifier
from sklearn.naive_bayes import MultinomialNB
from sklearn.naive_bayes import BernoulliNB

# Load SuperLearner
from mlens.ensemble import SuperLearner

# create a list of base-models
def get_models():
    models = list()
    models.append(MultinomialNB())
    models.append(BernoulliNB())
    models.append(LogisticRegression(solver = 'liblinear', max_iter = 200))
    models.append(DecisionTreeClassifier())
    models.append(SVC(gamma = 'scale', probability = True))
    models.append(GaussianNB())
    models.append(KNeighborsClassifier())
    models.append(AdaBoostClassifier())
    models.append(BaggingClassifier())
    models.append(RandomForestClassifier())
    models.append(ExtraTreesClassifier())
    return models

# create the super learner
def get_super_learner(X):
    ensemble = SuperLearner(scorer = accuracy_score, folds = 5, shuffle = True, sample_size = len(X))

    # add base models
    models = get_models()
    ensemble.add(models)

    # add the meta model
    ensemble.add_meta(LogisticRegression(solver = 'liblinear'))
    return ensemble

# load press release data
train_data = pd.read_csv("../train_data.csv", index_col = 0)
train_labels = pd.read_csv("../train_labels.csv", index_col = 0)

X = train_data.values # Documents are already in random order
y = np.asarray([int(i) for i in train_labels.values])

# create the inputs and outputs

# split data
train, test, train_val, test_val = train_test_split(X, y, test_size = 0.20) # test_size = 0.50
print('Train', train.shape, train_val, 'Test', test.shape, test_val.shape)

# create the super learner
ensemble = get_super_learner(train)

```

```

# fit the super learner
ensemble.fit(train, train_val)

# summarize base learners
print(ensemble.data)

# make predictions on hold out set
test_hat = ensemble.predict(test)

print('Super Learner: %.3f' % (accuracy_score(test_val, test_hat) * 100))

## [MLENS] backend: threading
## Train (2193, 3386) [ 7  3 20 ...  2 16 15] Test (549, 3386) (549,)
##
##          score-m  score-s   ft-m   ft-s   pt-m  pt-s
## layer-1 adaboostclassifier    0.19    0.01  37.82   3.02   0.80  0.12
## layer-1 baggingclassifier     0.50    0.02  48.02   4.01   0.62  0.22
## layer-1 bernoullinb          0.52    0.02   5.75   0.25   0.06  0.01
## layer-1 decisiontreeclassifier 0.38    0.02  31.90   1.62   0.01  0.01
## layer-1 extratreesclassifier    0.61    0.02  82.02   5.16   0.35  0.06
## layer-1 gaussiannb            0.39    0.03   0.74   0.22   1.22  0.23
## layer-1 kneighborsclassifier    0.10    0.01   0.01   0.02   0.86  0.18
## layer-1 logisticregression     0.61    0.01   6.34   1.96   0.03  0.02
## layer-1 multinomialnb         0.63    0.01  10.60   0.42   0.05  0.03
## layer-1 randomforestclassifier 0.58    0.01  22.91   1.66   0.32  0.07
## layer-1 svc                   0.51    0.03 397.54 119.21 11.45  4.96
##
## Super Learner: 25.501

```

The SuperLearner has a much lower accuracy than the base learners?!

4 Classification of unlabelled data

4.1 Using the NB textmodel

We trained the models using a set of 2,742 labelled documents. In order to obtain aggregated measures of issue attention, we predict the issue categories of all 47,111 labelled and unlabelled press releases in our sample.

```

tmod_nb_r1 <- textmodel_nb(dfmat_training, dfmat_training$issue_r1, distribution = "multinomial")

dfmat_matched <- dfm_match(dfmat_test,
                           features = featnames(dfmat_training))

actual_class <- docvars(dfmat_matched, "issue_r1")
predicted_class <- predict(tmod_nb_r1, newdata = dfmat_matched)
tab_class <- table(actual_class, predicted_class)
tab_class

```

```

##          predicted_class
## actual_class  1  2  3  4  5  6  7  9 10 12 15 16 17 20 99 191 192
##          1   24  0  0  0  4  1  4  0  1  1  2  0  1  0  0  0  1
##          2    0 20  2  1  0  1  0  2  0  7  0  2  0  2  1  3  0
##          3    1  1 16  0  0  1  1  0  0  0  0  0  0  0  0  1  0
##          4    0  0  0 16  0  0  1  0  0  0  0  0  0  0  0  0  0
##          5    4  0  2  0 23  0  0  1  0  1  2  0  0  0  0  0  1

```

```
##      6      0      1      0      0      2 20      0      0      3      0      0      0      0      0      1      0      0
##      7      0      0      0      4      0      0 29      0      3      1      1      0      1      1      0      1      1
##      9      0      0      0      0      1      1      0 13      0      1      0      0      0      2      1      1      0
##     10      1      1      0      1      3      1      5      1 17      1      6      0      0      1      1      1      0
##     12      0      2      0      1      0      1      2      3      0 21      2      1      1      1      1      2      2
##     15      1      3      1      0      0      1      3      0      2      2 13      0      1      0      1      1      9
##     16      0      2      1      0      0      0      0      0      0      2      1 15      0      2      0      3      0
##     17      0      2      0      0      0      3      0      0      0      1      0      0      4      0      2      0      0
##     20      3      0      0      0      0      2      1      0      0      1      0      1      1      3      0      0      0
##     99      2      0      0      0      0      4      0      0      0      0      0      2      1      2 20      2      1
##    191      0      1      0      1      0      1      3      0      0      1      0      5      0      3      0      55      4
##    192      2      0      0      0      0      0      0      1      0      1      2      0      0      0      2      3      20
```

```
# Loading full dataset from parent dir
all_germany <- read_rds("../data_joint.RDS") %>% select(c(header, text.x, date.x, issue, party.x, id))
```

```
# Constructing the document frequency matrix
dfmat_all <- corpus(str_c(all_germany$header, " ", all_germany$text.x)) %>%
  dfm(remove = stopwords("de"), # Stem and remove stopwords, punctuation etc.
      stem = T,
      remove_punct = T,
      remove_number = T,
      remove_symbols = T,
      remove_url = T)
```

```
## Warning: 'dfm.corpus()' is deprecated. Use 'tokens()' first.
```

```
## Warning: '...' should not be used for tokens() arguments; use 'tokens()' first.
```

```
## Warning: remove_number argument is not used.
```

```
## Warning: 'remove' is deprecated; use dfm_remove() instead
```

```
## Warning: 'stem' is deprecated; use dfm_wordstem() instead
```

```
# Adding docvars
docvars(dfmat_all, "party") <- all_germany$party.x
docvars(dfmat_all, "date") <- all_germany$date.x
docvars(dfmat_all, "id") <- all_germany$id
```

```
# Subsetting to features in the training data
dfmat_all <- dfm_match(dfmat_all, features = featnames(dfmat_training))
```

```
# Predicting the issue category for all documents
dfmat_all$issue_r1 <- predict(tmod_nb_r1, newdata = dfmat_all)
```

```
table(dfmat_all$issue_r1)
```

```
##
##      1      2      3      4      5      6      7      9     10     12     15     16     17     20     99    191
## 2714 2672 1703 1925 2925 3306 3361 2411 2997 3574 3086 1988 1122 1619 2652 6083
##    192
## 2973
```

4.2 Aggregation of the issues categories over time and party

To measure parties' evolving issue agendas, we aggregate the category counts over time.


```

# Create dataframe from dfm
issue_agendas <- data.frame(date = docvars(dfmat_all, "date"), party = docvars(dfmat_all,
  "party"), issue_r1 = docvars(dfmat_all, "issue_r1"))

# Make date quarterly
issue_agendas$date <- as.character(issue_agendas$date) %>%
  substr(1, 8) %>%
  str_c("15") %>%
  str_replace_all(c(`^-01-` = "-02-", `^-03-` = "-02-", `^-04-` = "-05-", `^-06-` = "-05-",
    `^-07-` = "-08-", `^-09-` = "-08-", `^-10-` = "-11-", `^-12-` = "-11-")) %>%
  ymd()

# Add variable for counting
issue_agendas$freq <- 1

# Aggregate by party, date and issue
issue_agendas <- aggregate(freq ~ party + date + issue_r1, issue_agendas, sum)

# Add var for total press releases per party and month
issue_agendas$party_sum <- ave(issue_agendas$freq, issue_agendas$date, issue_agendas$party,
  FUN = sum)

issue_agendas$attention <- issue_agendas$freq/issue_agendas$party_sum

# Add issue descriptions
issue_agendas <- merge(issue_agendas, issue_categories, by = "issue_r1")

```

4.3 Plotting issue attention: Examples

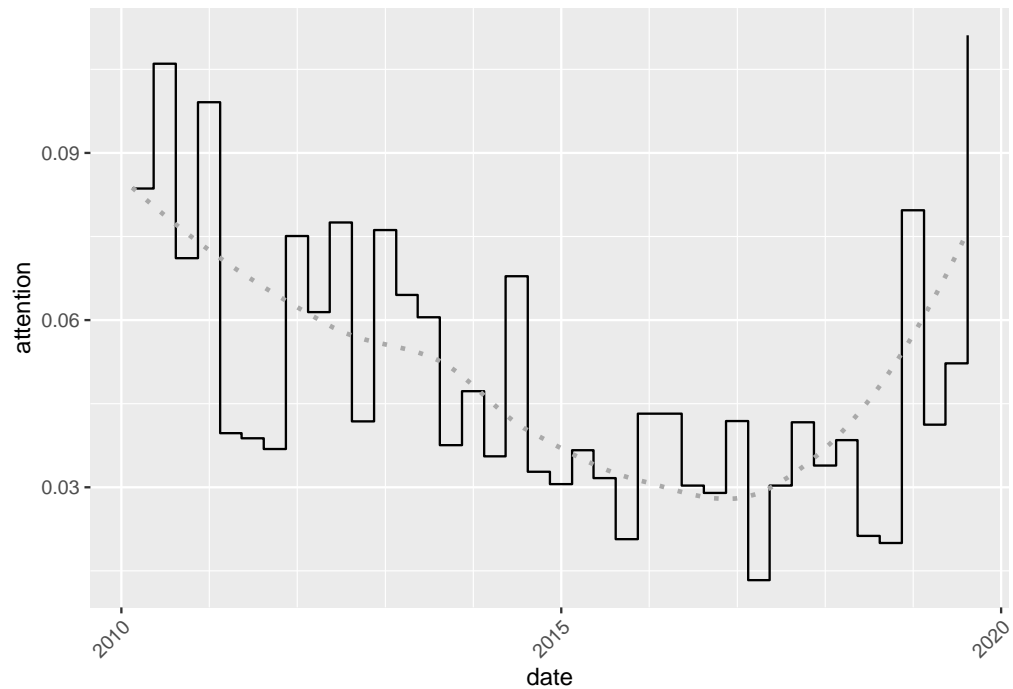
```

if (!dir.exists("plots")) dir.create("plots")

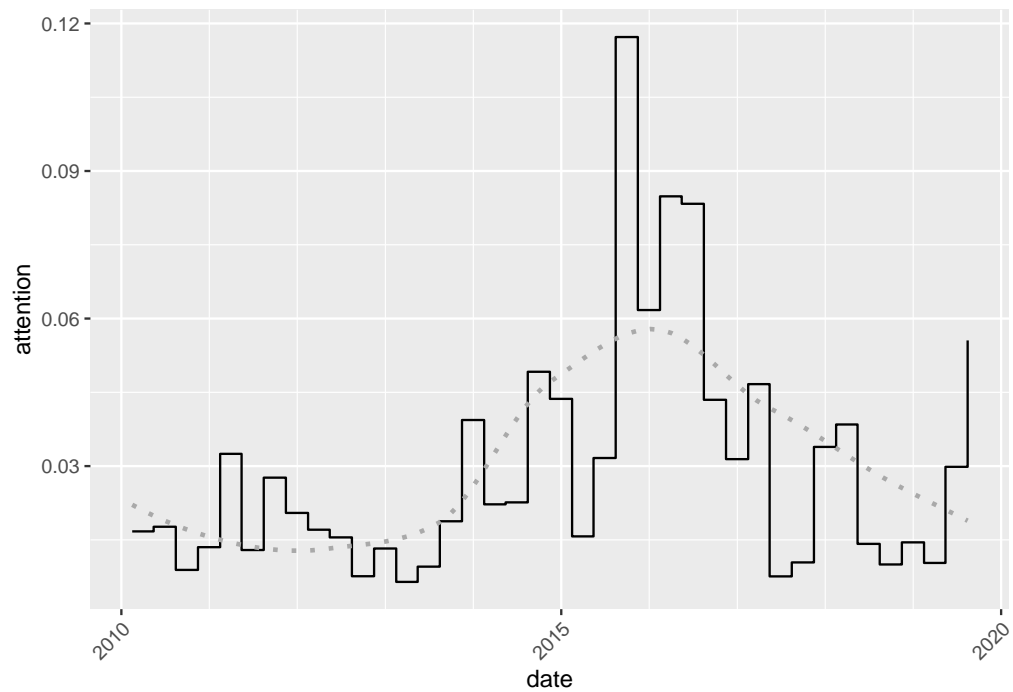
# Function for plotting parties' issue attention over time
plot_issue_party <- function(plot_issue, plot_party) ggplot(issue_agendas %>%
  filter(issue_r1 == plot_issue & party == plot_party), aes(x = date, y = attention)) +
  geom_step() + geom_smooth(method = "loess", formula = "y ~ x", color = "dark grey",
    lty = 3, se = F) + theme(axis.text.x = element_text(angle = 45, hjust = 1)) +
  scale_x_date(date_minor_breaks = "1 year") + # ggtitle('Share of press releases for issue per quart
# issue_categories$issue_r1_descr[issue_categories$issue_r1 == plot_issue], ' -
# ', plot_party, ')') +
ggsave(str_c("plots/", plot_issue, " - ", issue_categories$issue_r1_descr[issue_categories$issue_r1 ==
  plot_issue], "_", plot_party, ".pdf"), device = cairo_pdf, width = 5 * 2^0.5,
  height = 5) + ggsave(str_c("plots/", plot_issue, " - ", issue_categories$issue_r1_descr[issue_cate
  plot_issue], "_", plot_party, ".png"), width = 5 * 2^0.5, height = 5)

# Plot quarterly issue attention for category '7 Environment & Energy' for
# 'union_fraktion'
plot_issue_party(7, "union_fraktion") # There seems to be a decline since Fukushima in 2011.

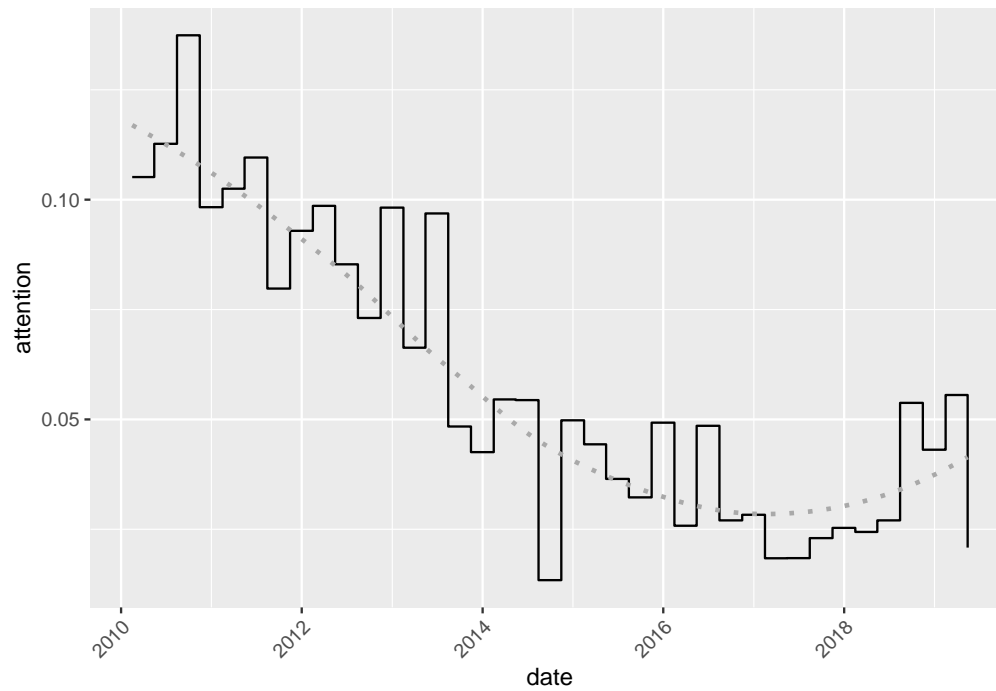
```



```
# Plot quarterly issue attention for category '9 Immigration' for
# 'union_fraktion'
plot_issue_party(9, "union_fraktion") # There seems to be a peak around the so-called refugee crisis i
```



```
# Plot quarterly issue attention for category '7 Environment & Energy' for
# 'spd_fraktion'
plot_issue_party(7, "spd_fraktion") # There seems to be a decline since Fukushima in 2011.
```



```
# Plot quarterly issue attention for category '10 Welfare' for 'spd_fraktion'
plot_issue_party(10, "spd_fraktion") # There seems to be a lower emphasis on welfare after the entry i
```

