

Feature selection

```
# library(extrafont)
# extrafont::loadfonts(quiet = TRUE)

set.seed(42)
library(igraph)

##
## Attaching package: 'igraph'

## The following objects are masked from 'package:stats':
##
##      decompose, spectrum

## The following object is masked from 'package:base':
##
##      union

library(tidyverse)

## -- Attaching core tidyverse packages ----- tidyverse 2.0.0 --
## v dplyr      1.1.4      v readr      2.1.5
## v forcats    1.0.0      v stringr   1.5.1
## v ggplot2     3.5.1      v tibble     3.2.1
## v lubridate   1.9.3      v tidyr      1.3.1
## v purrr       1.0.2

## -- Conflicts ----- tidyverse_conflicts() --
## x lubridate::%--%()      masks igraph::%--%()
## x dplyr::as_data_frame() masks tibble::as_data_frame(), igraph::as_data_frame()
## x purrr::compose()       masks igraph::compose()
## x tidyr::crossing()       masks igraph::crossing()
## x dplyr::filter()         masks stats::filter()
## x dplyr::lag()            masks stats::lag()
## x purrr::simplify()       masks igraph::simplify()
## i Use the conflicted package (<http://conflicted.r-lib.org/>) to force all conflicts to become errors
```

Load and tidy data

```
data <- read_csv("../measurements/measurements.csv")

## Rows: 754 Columns: 96
## -- Column specification -----
## Delimiter: ","
## chr  (9): fpath, KUK_ID, class, FileName, FolderPath, subcorpus, DocumentTit...
## dbl  (85): RuleAbstractNouns, RuleAmbiguousRegards, RuleAnaphoricReferences, ...
## lgl   (2): ClarityPursuit, SyllogismBased
##
## i Use `spec()` to retrieve the full column specification for this data.
```

i Specify the column types or set `show_col_types = FALSE` to quiet this message.

```
data_clean <- data %>%
  select(!c(
    fpath,
    # KUK_ID,
    # FileName,
    FolderPath,
    # subcorpus,
    DocumentTitle,
    ClarityPursuit,
    Readability,
    SyllogismBased,
    SourceDB
  )) %>%
  # replace -1s in variation coefficients with NAs
  mutate(across(c(
    `RuleDoubleAdpos.max_allowable_distance.v`,
    `RuleTooManyNegations.max_negation_frac.v`,
    `RuleTooManyNegations.max_allowable_negations.v`,
    `RuleTooManyNominalConstructions.max_noun_frac.v`,
    `RuleTooManyNominalConstructions.max_allowable_nouns.v`,
    `RuleCaseRepetition.max_repetition_count.v`,
    `RuleCaseRepetition.max_repetition_frac.v`,
    `RulePredSubjDistance.max_distance.v`,
    `RulePredObjDistance.max_distance.v`,
    `RuleInfVerbDistance.max_distance.v`,
    `RuleMultiPartVerbs.max_distance.v`,
    `RuleLongSentences.max_length.v`,
    `RulePredAtClauseBeginning.max_order.v`,
    `mattr.v`,
    `maentropy.v`
  ), ~ na_if(.x, -1))) %>%
  # replace NAs with 0s
  replace_na(list(
    RuleGPcoordovs = 0,
    RuleGPdeverbaddr = 0,
    RuleGPpatinstr = 0,
    RuleGPdeverbsubj = 0,
    RuleGPadjective = 0,
    RuleGPpatbenperson = 0,
    RuleGPwordorder = 0,
    RuleDoubleAdpos = 0,
    RuleDoubleAdpos.max_allowable_distance = 0,
    RuleDoubleAdpos.max_allowable_distance.v = 0,
    RuleAmbiguousRegards = 0,
    RuleReflexivePassWithAnimSubj = 0,
    RuleTooManyNegations = 0,
    RuleTooManyNegations.max_negation_frac = 0,
    RuleTooManyNegations.max_negation_frac.v = 0,
    RuleTooManyNegations.max_allowable_negations = 0,
    RuleTooManyNegations.max_allowable_negations.v = 0,
    RuleTooManyNominalConstructions.max_noun_frac.v = 0,
    RuleTooManyNominalConstructions.max_allowable_nouns.v = 0,
```

```

RuleFunctionWordRepetition = 0,
RuleCaseRepetition.max_repetition_count.v = 0,
RuleCaseRepetition.max_repetition_frac.v = 0,
RuleWeakMeaningWords = 0,
RuleAbstractNouns = 0,
RuleRelativisticExpressions = 0,
RuleConfirmationExpressions = 0,
RuleRedundantExpressions = 0,
RuleTooLongExpressions = 0,
RuleAnaphoricReferences = 0,
RuleLiteraryStyle = 0,
RulePassive = 0,
RulePredSubjDistance = 0,
RulePredSubjDistance.max_distance = 0,
RulePredSubjDistance.max_distance.v = 0,
RulePredObjDistance = 0,
RulePredObjDistance.max_distance = 0,
RulePredObjDistance.max_distance.v = 0,
RuleInfVerbDistance = 0,
RuleInfVerbDistance.max_distance = 0,
RuleInfVerbDistance.max_distance.v = 0,
RuleMultiPartVerbs = 0,
RuleMultiPartVerbs.max_distance = 0,
RuleMultiPartVerbs.max_distance.v = 0,
RuleLongSentences.max_length.v = 0,
RulePredAtClauseBeginning.max_order.v = 0,
RuleVerbalNouns = 0,
RuleDoubleComparison = 0,
RuleWrongValencyCase = 0,
RuleWrongVerbominalCase = 0,
RuleIncompleteConjunction = 0
)) %>%
# norm data expected to correlate with text length
mutate(across(c(
  RuleGPcoordovs,
  RuleGPdeverbaddr,
  RuleGPpatinstr,
  RuleGPdeverbsubj,
  RuleGPadjective,
  RuleGPpatbenperson,
  RuleGPwordorder,
  RuleDoubleAdpos,
  RuleAmbiguousRegards,
  RuleFunctionWordRepetition,
  RuleWeakMeaningWords,
  RuleAbstractNouns,
  RuleRelativisticExpressions,
  RuleConfirmationExpressions,
  RuleRedundantExpressions,
  RuleTooLongExpressions,
  RuleAnaphoricReferences,
  RuleLiteraryStyle,
  RulePassive,

```

```

RuleVerbalNouns,
RuleDoubleComparison,
RuleWrongValencyCase,
RuleWrongVerbNominalCase,
RuleIncompleteConjunction,
num_hapax,
RuleReflexivePassWithAnimSubj,
RuleTooManyNominalConstructions,
RulePredSubjDistance,
RuleMultiPartVerbs,
RulePredAtClauseBeginning
), ~ .x / word_count)) %>%
mutate(across(c(
  RuleTooFewVerbs,
  RuleTooManyNegations,
  RuleCaseRepetition,
  RuleLongSentences,
  RulePredObjDistance,
  RuleInfVerbDistance
), ~ .x / sent_count)) %>%
# remove variables identified as "u counts"
select(!c(
  RuleTooFewVerbs,
  RuleTooManyNegations,
  RuleTooManyNominalConstructions,
  RuleCaseRepetition,
  RuleLongSentences,
  RulePredAtClauseBeginning,
  sent_count,
  word_count,
  syllab_count,
  char_count
)) %>%
# remove variables identified as unreliable
select(!c(
  RuleAmbiguousRegards,
  RuleFunctionWordRepetition,
  RuleDoubleComparison,
  RuleWrongValencyCase,
  RuleWrongVerbNominalCase
)) %>%
# remove further variables belonging to the 'acceptability' category
select(!RuleIncompleteConjunction) %>%
mutate(across(c(class), ~ as.factor(.x)))

# no NAs should be present now
data_clean[!complete.cases(data_clean), ]

## # A tibble: 0 x 73
## # i 73 variables: KUK_ID <chr>, class <fct>, FileName <chr>, subcorpus <chr>,
## #   RuleAbstractNouns <dbl>, RuleAnaphoricReferences <dbl>,
## #   RuleCaseRepetition.max_repetition_count <dbl>,
## #   RuleCaseRepetition.max_repetition_count.v <dbl>,
## #   RuleCaseRepetition.max_repetition_frac <dbl>,

```

```
## # RuleCaseRepetition.max_repetition_frac.v <dbl>,
## # RuleConfirmationExpressions <dbl>, RuleDoubleAdpos <dbl>, ...

data_clean_scaled <- data_clean %>%
  mutate(across(class, ~ .x == "good")) %>%
  mutate(across(5:length(names(data_clean)), ~ scale(.x)))
```

Important features identification

```
data_clean_good <- data_clean_scaled %>% filter(class == "good")
data_clean_bad <- data_clean_scaled %>% filter(class == "bad")

feature_importances <- tibble(
  feat_name = character(), p_value = numeric()
)

for (i in 5:73) {
  fname <- names(data_clean)[i]

  formula_single <- reformulate(fname, "class")
  # print(formula_single)

  glm_model <- glm(formula_single, data_clean, family = "binomial")
  glm_coefficients <- summary(glm_model)$coefficients
  row_index <- which(rownames(glm_coefficients) == fname)
  p_value <- glm_coefficients[row_index, 4]

  feature_importances <- feature_importances %>%
    add_row(feat_name = fname, p_value = p_value)
}
feature_importances
```

```
## # A tibble: 69 x 2
##   feat_name                                p_value
##   <chr>                                <dbl>
## 1 RuleAbstractNouns                      0.00187
## 2 RuleAnaphoricReferences                0.660
## 3 RuleCaseRepetition.max_repetition_count 0.0722
## 4 RuleCaseRepetition.max_repetition_count.v 0.00479
## 5 RuleCaseRepetition.max_repetition_frac  0.000000740
## 6 RuleCaseRepetition.max_repetition_frac.v 0.000000472
## 7 RuleConfirmationExpressions            0.0985
## 8 RuleDoubleAdpos                       0.312
## 9 RuleDoubleAdpos.max_allowable_distance 0.000154
## 10 RuleDoubleAdpos.max_allowable_distance.v 0.00000356
## # i 59 more rows
```

```
selected_features <- feature_importances %>%
  filter(p_value <= 0.05) %>%
  pull(feat_name)
```

Correlations

```
data_pure <- data_clean %>% select(any_of(selected_features))

cor_matrix <- cor(data_pure)

cor_tibble_long <- cor_matrix %>%
  as_tibble() %>%
  mutate(feat1 = selected_features) %>%
  pivot_longer(!feat1, names_to = "feat2", values_to = "cor") %>%
  mutate(abs_cor = abs(cor))

cor_matrix_upper <- cor_matrix
cor_matrix_upper[lower.tri(cor_matrix_upper)] <- 0

cor_tibble_long_upper <- cor_matrix_upper %>%
  as_tibble() %>%
  mutate(feat1 = selected_features) %>%
  pivot_longer(!feat1, names_to = "feat2", values_to = "cor") %>%
  mutate(abs_cor = abs(cor)) %>%
  filter(feat1 != feat2)
cor_tibble_long_upper_considerable <- cor_tibble_long_upper %>%
  filter(abs_cor >= 0.3)

max_correlations <- cor_tibble_long %>%
  filter(feat1 != feat2) %>%
  group_by(feat1) %>%
  summarize(maxcor = max(abs_cor)) %>%
  ungroup()
```

Visualisation

```
library(paletteer)

my_colors <- paletteer::paletteer_d("ggthemes::Classic_10_Medium")

network <- graph_from_data_frame(
  cor_tibble_long_upper_considerable,
  directed = FALSE
)
E(network)$weight <- cor_tibble_long_upper_considerable$abs_cor
network_communities <- cluster_optimal(network)

network_membership <- membership(network_communities)

plot(
  network,
  layout = layout_fruchterman_reingold,
  vertex_color = map(
    network_communities$membership,
    function(x) my_colors[x]
  ) %>% unlist(use.names = FALSE),
  vertex_size = 6,
```

```

# vertex.frame.color = "#00000000",
# vertex.label.family = "Public Sans",
vertex.label.color = "black",
vertex.label.cex = 0.5
)

```

