

Oppgave 3

Formål

Formålet med denne oppgaven er å belyse i hvilke, og i hvilken grad, partiene på stortinget endret sine posisjoner mellom 2001 og 2005. Du skal først laste ned partiprogrammene for alle partiene på stortinget i 2001 og 2005. Du finner alle partiprogrammene her: Du finner en oversik over partier på stortinget her:

Last inn partiprogrammene

I snutten nedenfor skal du skrive en kode som laster inn alle partiprogrammene, gjør dem om til et corpus, legger til partinavn og årstall som dokument variabler, og lagrer dette som et r-object som heter `valgprogram.RData`. Bruk funksjoner i R-pakkene `readtext` og `quanteda`. Når du har gjort denne riktig, setter du `eval = FALSE` i topplinjen på `last_ned` biten med R-kode slik at du ikke må laste inn data hver gang du kjører koden. Du kan sette tilbake til `eval = TRUE` før du leverer inn.

```
# 1 poeng
url <- c(
  "http://www.nsd.uib.no/polsys/data/filer/parti/H3.html", # Høyre 2001
  "http://www.nsd.uib.no/polsys/data/filer/parti/H19.html", # AP 2001
  "http://www.nsd.uib.no/polsys/data/filer/parti/H25.html", # Venstre 2001
  "http://www.nsd.uib.no/polsys/data/filer/parti/H32.html", # SV 2001
  "http://www.nsd.uib.no/polsys/data/filer/parti/H34.html", # KrF 2001
  "http://www.nsd.uib.no/polsys/data/filer/parti/H37.html", # Senterpartiet 2001
  "http://www.nsd.uib.no/polsys/data/filer/parti/H5181.html", # FrP 2001
  "http://www.nsd.uib.no/polsys/data/filer/parti/H9364.html", # Høyre 2005
  "http://www.nsd.uib.no/polsys/data/filer/parti/H9191.html", # Ap 2005
  "http://www.nsd.uib.no/polsys/data/filer/parti/H9362.html", # Venstre 2005
  "http://www.nsd.uib.no/polsys/data/filer/parti/H9360.html", # SV 2005
  "http://www.nsd.uib.no/polsys/data/filer/parti/H9365.html", # KrF 2005
  "http://www.nsd.uib.no/polsys/data/filer/parti/H9366.html", # Senterpartiet 2005
  "http://www.nsd.uib.no/polsys/data/filer/parti/H9368.html" # FrP 2005
)
partinavn <- c("H","Ap","V","SV","KrF","Sp","FrP","H","Ap","V","SV","KrF","Sp","FrP")
year <- c(rep(2001,7),rep(2005,7))
url %>%
  readtext(encoding = "LATIN1") %>%
  corpus() -> valgprogram
docvars(valgprogram, "parti") <- partinavn
docvars(valgprogram, "year") <- year
docvars(valgprogram, "navn") <- paste(partinavn,year, sep = "_")
save(valgprogram, file = "valgprogram.RData")
```

Lengde på partiprogrammer

Nå vil jeg at du skal lage en tabell som viser hvor lange de ulike partiprogrammene er, hvor mange ulike ord partiene bruker, og forholdet mellom antall ord og antall setninger i parti-programmet. Bruk redskap fra `tidyverse` og `tidytext` for å gjøre dette. Du skal ikke gjøre noe preprocessing før du gjør dette.

```
## 3 poeng
```

```

load("valgprogram.RData")
summary(valgprogram)

## Corpus consisting of 14 documents:
##
##      Text Types Tokens Sentences parti year   navn
##   H34.html  7331  47109     2776    H 2001   H_2001
## H9362.html  7565  47566     2563   Ap 2001  Ap_2001
## H9191.html  4752  25962     1162    V 2001  V_2001
## H9368.html  7787  48604     3019   SV 2001  SV_2001
## H9365.html  5804  33181     1505  KrF 2001  KrF_2001
## H9366.html  9586  66535     2784    Sp 2001  Sp_2001
##   H32.html  6878  37929     1458  FrP 2001  FrP_2001
##   H25.html  7071  44327     2336    H 2005   H_2005
## H9364.html  4448  24844     1433   Ap 2005  Ap_2005
##   H37.html  8357  53093     2124    V 2005  V_2005
## H9360.html  6127  36678     1972   SV 2005  SV_2005
##   H19.html  6567  43333     2446  KrF 2005  KrF_2005
##    H3.html  4175  20473     1176    Sp 2005  Sp_2005
## H5181.html  7166  41343     2146  FrP 2005  FrP_2005
##
## Source: /Users/bjornkho/teaching/stv1515/v19/sensur/oppgabe_3/* on x86_64 by bjornkho
## Created: Thu Apr 18 12:04:27 2019
## Notes:
## tell opp antall ord per parti
dfm(valgprogram) %>%
  tidy() %>%
  rename("n" = count) -> tmp
tmp %>%
  group_by(document) %>%
  summarise(lengde = sum(n)) %>%
  left_join(tmp) -> parti_ord

## Joining, by = "document"
# tell opp antall setninger per parti
docnames(valgprogram)

## [1] "H34.html" "H9362.html" "H9191.html" "H9368.html" "H9365.html"
## [6] "H9366.html" "H32.html" "H25.html" "H9364.html" "H37.html"
## [11] "H9360.html" "H19.html" "H3.html" "H5181.html"

out <- rep(NA, length(docnames(valgprogram)))
for (i in 1:length(docnames(valgprogram))){
temp <- data_frame(txt = texts(valgprogram[docnames(valgprogram)[i]]))
out[i] <- temp %>%
  unnest_tokens("setninger", txt, token = "sentences") %>%
  nrow()
}

## Warning: `data_frame()` is deprecated, use `tibble()`.
## This warning is displayed once per session.
# sett sammen of lag informative partinavn og år
tibble(document = docnames(valgprogram), n_setninger = out) %>%
  right_join(parti_ord) %>%

```

```
mutate(
  document = case_when(
    document == "H3.html" ~ "H_01",
    document == "H19.html" ~ "Ap_01",
    document == "H25.html" ~ "V_01",
    document == "H32.html" ~ "SV_01",
    document == "H34.html" ~ "KrF_01",
    document == "H37.html" ~ "Sp_01",
    document == "H5181.html" ~ "FrP_01",
    document == "H9364.html" ~ "H_05",
    document == "H9191.html" ~ "Ap_05",
    document == "H9362.html" ~ "V_05",
    document == "H9360.html" ~ "SV_05",
    document == "H9365.html" ~ "KrF_05",
    document == "H9366.html" ~ "Sp_05",
    document == "H9368.html" ~ "FrP_05"),
) -> parti_ord
```

```
## Joining, by = "document"
```

```
## sett sammen til en oppsummerende tabell
```

```
parti_ord %>%
  group_by(document) %>%
  summarise(n_setninger = unique(n_setninger),
            n_ord = unique(lengde),
            snitt_setning = n_ord/n_setninger) %>%
  ungroup() %>%
  arrange(desc(snitt_setning)) -> oppsummering_parti
oppsummering_parti
```

```
## # A tibble: 14 x 4
##   document n_setninger n_ord snitt_setning
##   <chr>      <int> <dbl>      <dbl>
## 1 SV_01         1458 37929         26.0
## 2 Sp_01         2124 53093         25.0
## 3 Sp_05         2784 66535         23.9
## 4 Ap_05         1162 25962         22.3
## 5 KrF_05        1505 33181         22.0
## 6 FrP_01        2146 41343         19.3
## 7 V_01          2336 44327         19.0
## 8 SV_05         1972 36678         18.6
## 9 V_05          2563 47566         18.6
## 10 Ap_01         2446 43333         17.7
## 11 H_01          1176 20473         17.4
## 12 H_05          1433 24844         17.3
## 13 KrF_01        2776 47109         17.0
## 14 FrP_05        3019 48604         16.1
```

Fordeling av ordfrekvenser

Bruk ggplot for å lage et sett av figurer som viser hvordan ord-frekvenser fordeler seg per parti når begge valgprogrammene er sett samlet.

```

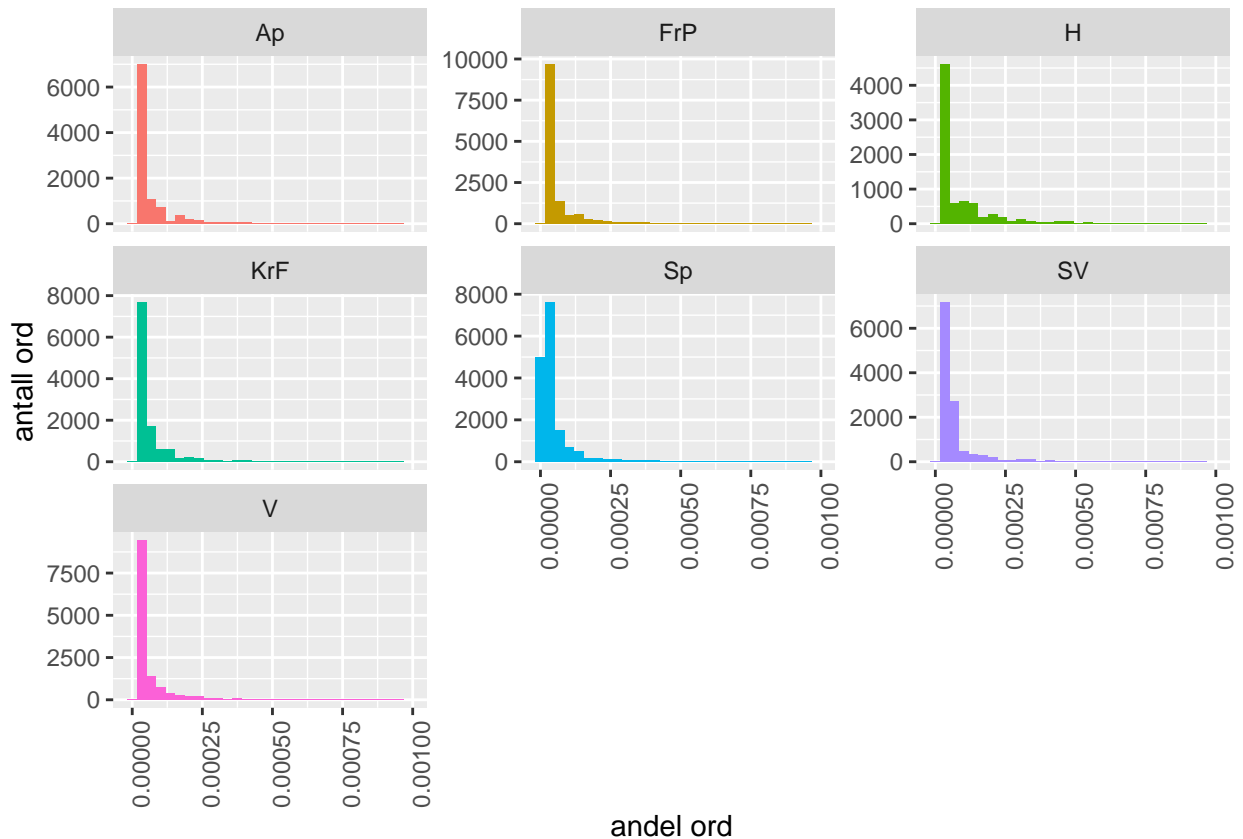
# 2 poeng
parti_ord %>%
  mutate(parti = case_when(
    document == "H_01" ~ "H",
    document == "Ap_01" ~ "Ap",
    document == "V_01" ~ "V",
    document == "SV_01" ~ "SV",
    document == "KrF_01" ~ "KrF",
    document == "Sp_01" ~ "Sp",
    document == "FrP_01" ~ "FrP",
    document == "H_05" ~ "H",
    document == "Ap_05" ~ "Ap",
    document == "V_05" ~ "V",
    document == "SV_05" ~ "SV",
    document == "KrF_05" ~ "KrF",
    document == "Sp_05" ~ "Sp",
    document == "FrP_05" ~ "FrP"),
  arstall = case_when(
    document == "H_01" ~ "2001",
    document == "Ap_01" ~ "2001",
    document == "V_01" ~ "2001",
    document == "SV_01" ~ "2001",
    document == "KrF_01" ~ "2001",
    document == "Sp_01" ~ "2001",
    document == "FrP_01" ~ "2001",
    document == "H_05" ~ "2005",
    document == "Ap_05" ~ "2005",
    document == "V_05" ~ "2005",
    document == "SV_05" ~ "2005",
    document == "KrF_05" ~ "2005",
    document == "Sp_05" ~ "2005",
    document == "FrP_05" ~ "2005")
  ) -> parti_ord
ggplot(parti_ord, aes(x=n/lengde, fill = parti)) +
  geom_histogram(show.legend = FALSE) +
  facet_wrap(~parti, ncol = 3, scales = "free_y") +
  theme(axis.text.x = element_text(angle = 90)) +
  scale_x_continuous(limits = c(NA,.001), name = "andel ord") +
  scale_y_continuous(name = "antall ord")

```

```

## `stat_bin()` using `bins = 30`. Pick better value with `binwidth`.
## Warning: Removed 1471 rows containing non-finite values (stat_bin).
## Warning: Removed 7 rows containing missing values (geom_bar).

```



Har senterpartiet endret seg?

Lag en figur som viser i hvilken grad Senterpartiet bruker andre ord i 2005 enn i 2001.

```
# 1 poeng
parti_ord %>%
  filter(document == c("Sp_01", "Sp_05")) %>% # kun Ap og Sp
  group_by(document) %>% # grupper ord per parti
  mutate(proportion = n / sum(n)) %>% #regn ut andel for hver av disse partiene for hvert ord
  select(- c(n_setninger, n, lengde, parti, arstall)) %>% # kast det vi ikke trenger
  spread(document, proportion) %>% # lag et bredt datasett
  gather(document, proportion, `Sp_05`) -> ord_prop # sett sammen til et langt datasett

library(scales)

##
## Attaching package: 'scales'
##
## The following object is masked from 'package:purrr':
##
##   discard
##
## The following object is masked from 'package:readr':
##
##   col_factor

ggplot(ord_prop, aes(x = proportion, y = `Sp_01`,
                    color = abs(`Sp_01` - proportion))) +
```

```

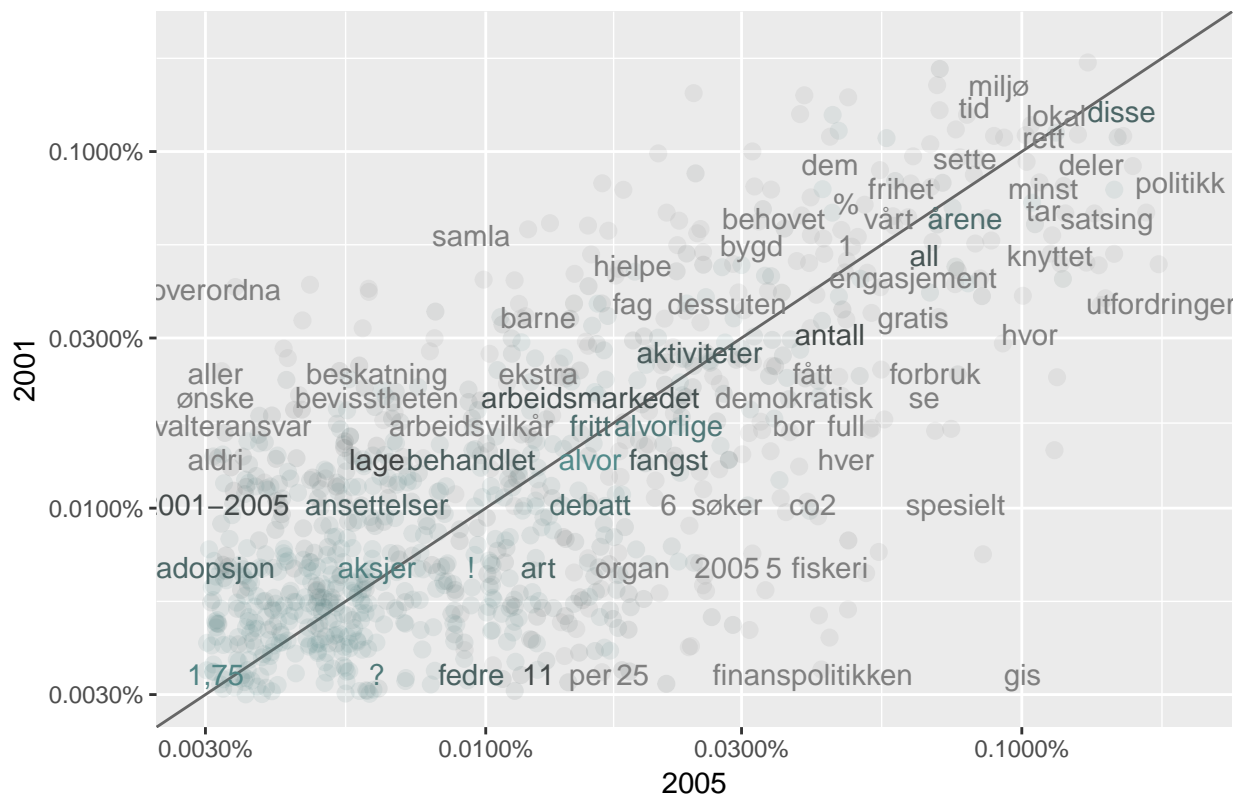
geom_abline(color = "grey40") + # 45 graders linje
geom_jitter(alpha = 0.1, size = 2.5, width = 0.3, height = 0.3) + # punkt per ord
geom_text(aes(label = term), check_overlap = TRUE, vjust = 1.5) + # legg på ord
scale_x_log10(labels = percent_format(), limits = c(0.00003,.002)) + # område for y-skalaen i log form
scale_y_log10(labels = percent_format(), limits = c(0.00003,.002)) + # område for x-skalaen i log form
scale_color_gradient(limits = c(0, 0.0001), # gradering av farer
                     low = "darkslategray4",
                     high = "grey25") +
theme(legend.position = "none") + # ingen nøkkel
labs(y = "2001", x = "2005") + # navn på y og x akse
ggtitle("Forskjell i ordbruk i SP fra 2001 til 2005")

```

Warning: Removed 6099 rows containing missing values (geom_point).

Warning: Removed 5742 rows containing missing values (geom_text).

Forskjell i ordbruk i SP fra 2001 til 2005



Hvilke ord skiller

Lag en dokument-frekvens matrise hvor du har tatt bort regnsetting, tall, symboler og urler. Gjør alt om til små bokstaver. Kalkuler tf og idf mål. Lag en figur som viser de 10 ordene som best skiller partiene fra hverandre.

```

# 2 poeng
valgprogram %>%
  tokens(what = "word",
         remove_numbers = TRUE,
         remove_punct = TRUE,

```

```

    remove_symbols = TRUE,
    remove_separators = TRUE,
    remove_url = TRUE) %>%
dfm() %>% # gjør om til dokument frekvens matrise
tidy() %>%
  rename("n" = count) %>%
  group_by(document) %>%
summarise(lengde = sum(n)) %>%
left_join(tmp) %>%
  mutate(
    document = case_when(
      document == "H3.html" ~ "H_01",
      document == "H19.html" ~ "Ap_01",
      document == "H25.html" ~ "V_01",
      document == "H32.html" ~ "SV_01",
      document == "H34.html" ~ "KrF_01",
      document == "H37.html" ~ "Sp_01",
      document == "H5181.html" ~ "FrP_01",
      document == "H9364.html" ~ "H_05",
      document == "H9191.html" ~ "Ap_05",
      document == "H9362.html" ~ "V_05",
      document == "H9360.html" ~ "SV_05",
      document == "H9365.html" ~ "KrF_05",
      document == "H9366.html" ~ "Sp_05",
      document == "H9368.html" ~ "FrP_05"),
    parti = case_when(
      document == "H_01" ~ "H",
      document == "Ap_01" ~ "Ap",
      document == "V_01" ~ "V",
      document == "SV_01" ~ "SV",
      document == "KrF_01" ~ "KrF",
      document == "Sp_01" ~ "Sp",
      document == "FrP_01" ~ "FrP",
      document == "H_05" ~ "H",
      document == "Ap_05" ~ "Ap",
      document == "V_05" ~ "V",
      document == "SV_05" ~ "SV",
      document == "KrF_05" ~ "KrF",
      document == "Sp_05" ~ "Sp",
      document == "FrP_05" ~ "FrP"),
    arstall = case_when(
      document == "H_01" ~ "2001",
      document == "Ap_01" ~ "2001",
      document == "V_01" ~ "2001",
      document == "SV_01" ~ "2001",
      document == "KrF_01" ~ "2001",
      document == "Sp_01" ~ "2001",
      document == "FrP_01" ~ "2001",
      document == "H_05" ~ "2005",
      document == "Ap_05" ~ "2005",
      document == "V_05" ~ "2005",
      document == "SV_05" ~ "2005",
      document == "KrF_05" ~ "2005",

```

```

document == "Sp_05" ~ "2005",
document == "FrP_05" ~ "2005")
) %>%
bind_tf_idf(term, document, n) -> parti_ren # lag tf_idf

```

Joining, by = "document"

```

parti_ren %>%
  arrange(desc(tf_idf)) %>% # sorter etter tf_idf
  mutate(term = factor(term, levels = rev(unique(term)))) %>% # lag en kategorisk variabel for hvert or
  group_by(parti) %>% # gruper etter parti
  top_n(10) %>% # veldg topp 8
  ungroup() -> ord1 # sett sammen til ett datasett, kall det ord1

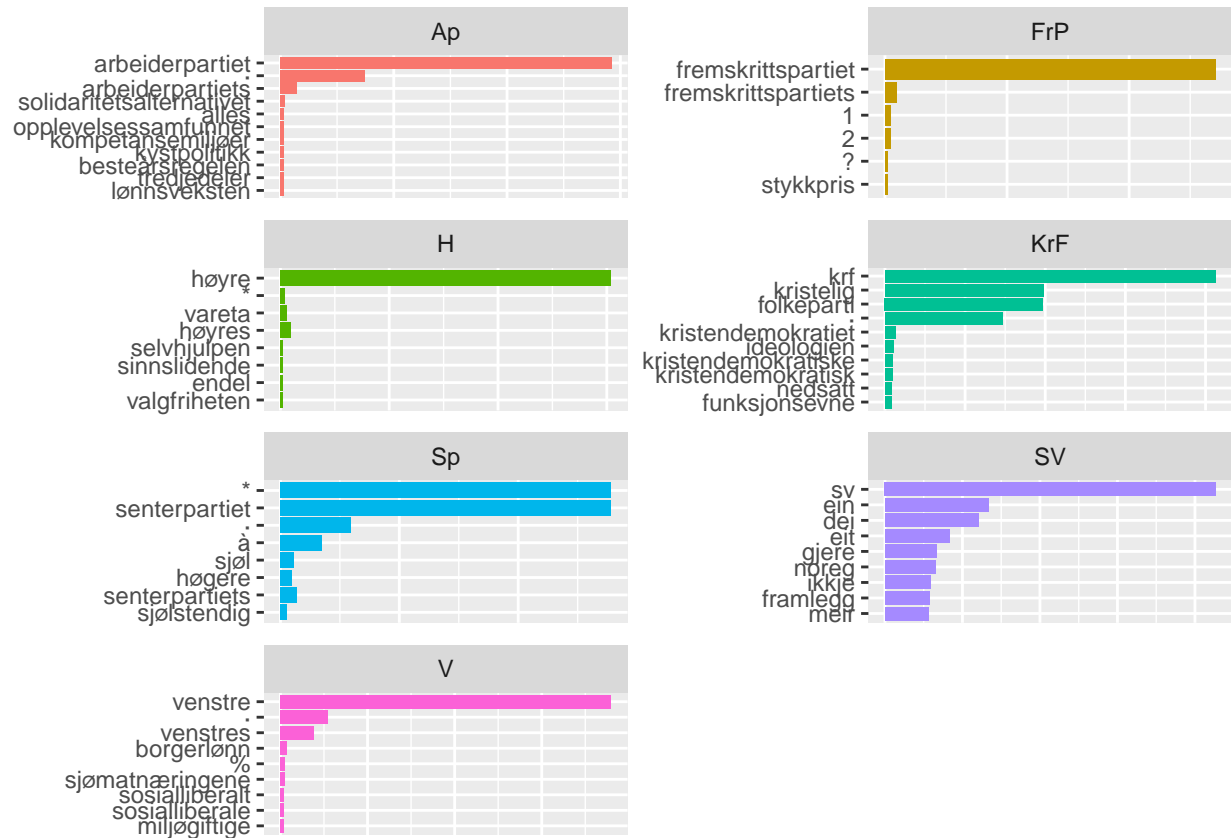
```

Selecting by tf_idf

```

ord1 %>%
  ggplot(aes(term, tf_idf, fill = parti)) + # lag figur
  geom_col(show.legend = FALSE) + # stolpediagram
  facet_wrap(~parti, ncol = 2, scale = "free") + # ett per parti
  theme(axis.title.x=element_blank(), # ta bort info å aksene
        axis.text.x=element_blank(),
        axis.ticks.x=element_blank(),
        axis.title.y = element_blank()) +
  coord_flip() -> fig_unfixed
fig_unfixed

```

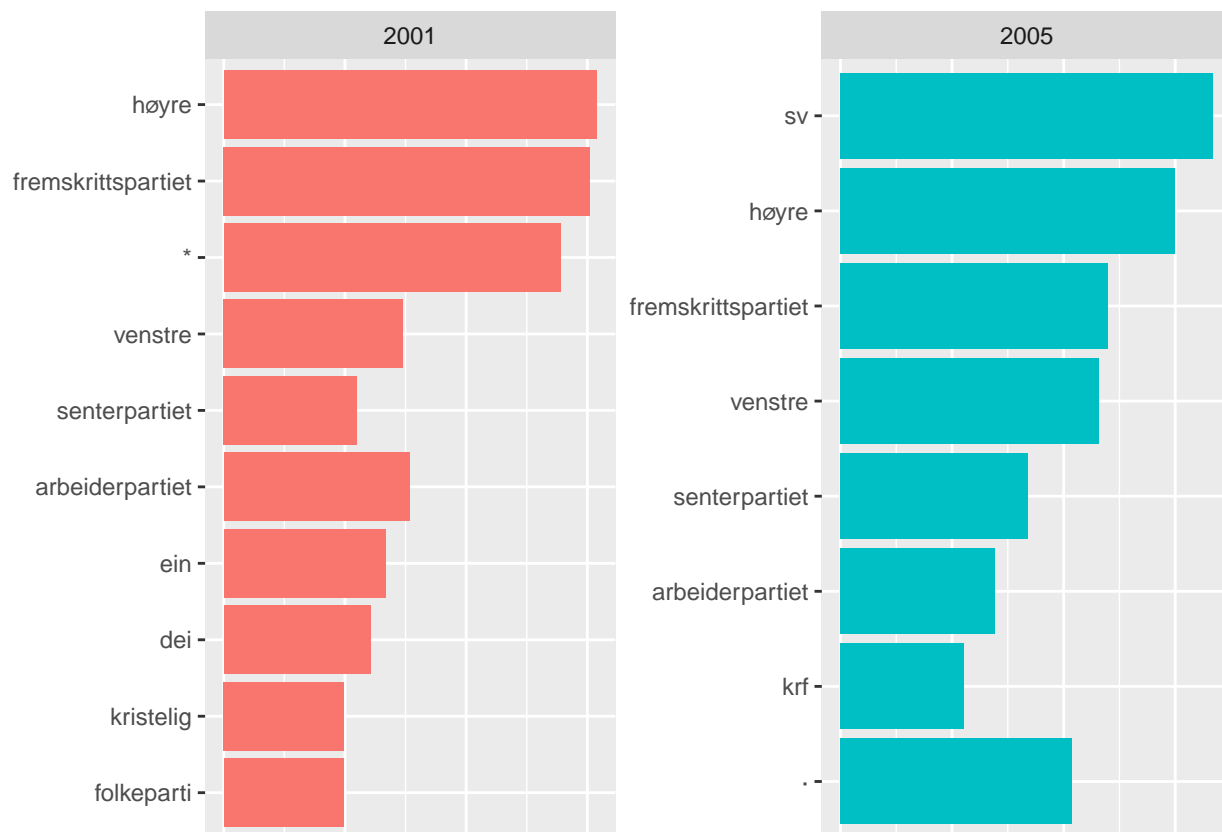


Gjennta det samme for årstall.


```
# 2 poeng
parti_ren %>%
  arrange(desc(tf_idf)) %>% # sorter etter tf_idf
  mutate(term = factor(term, levels = rev(unique(term)))) %>% # lag en kategorisk variabel for hvert or
  group_by(arstall) %>% # grupper etter parti
  top_n(10) %>% # veldg topp 8
  ungroup() -> ord1 # sett sammen til ett datasett, kall det ord1
```

```
## Selecting by tf_idf
```

```
ord1 %>%
  ggplot(aes(term, tf_idf, fill = arstall)) + # lag figur
  geom_col(show.legend = FALSE) + # stolpediagram
  facet_wrap(~arstall, scale = "free") + # ett per parti
  theme(axis.title.x=element_blank(), # ta bort info å aksene
        axis.text.x=element_blank(),
        axis.ticks.x=element_blank(),
        axis.title.y = element_blank()) +
  coord_flip()
```



Lag en variable med de ordene du mener er parti eller årstall-markører. Med markør menes et ord som i veldig stor grad identifiserer ett bestemt parti eller et årstall. Bruk denne variabelen til å fjerne disse ordene og lag en et nytt datasett uten disse ordene. Lag en figur som, per parti, viser topp 10 ord før og etter fjerning av disse markør-ordene.

```
# 2 poeng
# vi ser at SV sitt programm fra 2001 er på nynorsk. Vi dropper dette fra den videre analysen
parti_ren <- subset(parti_ren, drop(document != "SV_01"))
```

```

parti_ren %>%
  arrange(desc(tf_idf)) %>% # dorteer på tf_idf
  mutate(term = factor(term, levels = rev(unique(term)))) %>% # synkende
  group_by(parti) %>% # grupper etter parti
  top_n(25) %>% # de 25 ordene per parti som skiller mest
  pull(term) -> problem_ord# hent disse

```

```
## Selecting by tf_idf
```

```
problem_ord
```

```

## [1] sv høyre
## [3] fremskrittspartiet høyre
## [5] * fremskrittspartiet
## [7] venstre senterpartiet
## [9] arbeiderpartiet venstre
## [11] arbeiderpartiet krf
## [13] senterpartiet kristelig
## [15] folkeparti krf
## [17] . .
## [19] . à
## [21] . venstres
## [23] . svs
## [25] venstres vareta
## [27] sjøl fremskrittspartiets
## [29] høyres høgere
## [31] fremskrittspartiets senterpartiets
## [33] * arbeiderpartiets
## [35] høyres borgerlønn
## [37] kristendemokratiet foreslå
## [39] arbeiderpartiets ideologien
## [41] senterpartiets 1
## [43] 2 selvhjulpen
## [45] sinnslidende endel
## [47] sjølstendig %
## [49] sjømatnæringene samla
## [51] 1 kristendemokratiske
## [53] ? kristendemokratisk
## [55] valgfriheten sosialliberalt
## [57] sosialliberale miljøgiftige
## [59] stykkpris www.caplex.net
## [61] nedsatt funksjonsevne
## [63] arbeidstagere verdensorden
## [65] » 2
## [67] solidaritetsalternativet varetas
## [69] helsereformen «
## [71] » fagbevegelsen
## [73] krfs .
## [75] innovasjon trygghetsreform
## [77] " høyresidas
## [79] 4 kvinner
## [81] menneskesynet skattereform
## [83] overordna breiband
## [85] 6 privatisering

```

```

## [87] flat                sosialliberal
## [89] nedenfra              transportdirektoratet
## [91] læringsresultater    foreldrevalgte
## [93] 7                     ung
## [95] borgerne             høyesterett
## [97] premie               folkepensjonen
## [99] brei                 metodefrihet
## [101] utdannelses          tjenesteapparatet
## [103] brukervalg           alles
## [105] liberalt             språkteknologi
## [107] liberale             opplevelsessamfunnet
## [109] kompetansemiljøer   kystpolitikk
## [111] besteårsregelen      tredjedeler
## [113] lønnsveksten         holdbar
## [115] 2015                 helsehjelp
## [117] arbeidsmarkedsetaten velferden
## [119] sanksjonsmuligheter kystområdene
## [121] kristendemokratisk  nærhetsprinsippet
## [123] sjølstendige         4
## [125] "                    fins
## [127] avvikles             justeres
## [129] høyresida            privatiseringen
## [131] urett                varehandelen
## [133] særavgiftene         importerer
## [135] avbryte              atferdsproblemer
## [137] narkotikakriminalitet "
## [139] høyesterett          folketrygdfondets
## [141] genressurser         liberal
## [143] serviceerklæringer  gjenoppbygging
## [145] o                    6
## [147] livsform             tjenstepensjon
## [149] kapittel             tida
## [151] solidaritet          sjølstendighet
## [153] «                    »
## [155] minoritetsbakgrunn  menneskesynet
## [157] helsepersonellet    kapittel
## [159] foreldrepermisjonen kunstproduksjon
## [161] fiskerinæringa      forvalterskap
## [163] solidaritetsprinsippet pellets
## [165] forlengelse         samlivstiltak
## [167] idealisme           budsjettpolitikken
## [169] finmasket           satsingsområder
## [171] urettferdighet     ufrihet
## [173] aetat              tjenstepensjoner
## [175] levealder           nødhjelpsfond
## [177] rikdommer           nettoeksportør
## [179] arbeiderstyrte     kildekode
## [181] kulturnæringene    arbeidspress
## [183] fredsnasjon        åra
## [185] dumping             profitt
## 26387 Levels: langsiktige skolene arbeidsvilkår foregå lån ... sv
problem_ord <- tibble(
  term = as.character(problem_ord[-c(26,27,30,36,38,40,44:47,55,58,59,61:64,68,69,75,76,80:84,

```

```

86,87,89:92,94:106,108:114,116:120,122, 123,126:128,
130:137, 138:141, 143,144, 147:151, 152:155, 160:166, 168:189,
192:209, 210, 211, 213:215]))
parti_ren2 <- parti_ren%>%
  anti_join(problem_ord) # kast ut id ord

## Joining, by = "term"

parti_ren2 %>%
  arrange(desc(tf_idf)) %>%
  mutate(term = factor(term, levels = rev(unique(term)))) %>%
  group_by(parti) %>%
  top_n(9) %>%
  ungroup() -> ord2

## Selecting by tf_idf

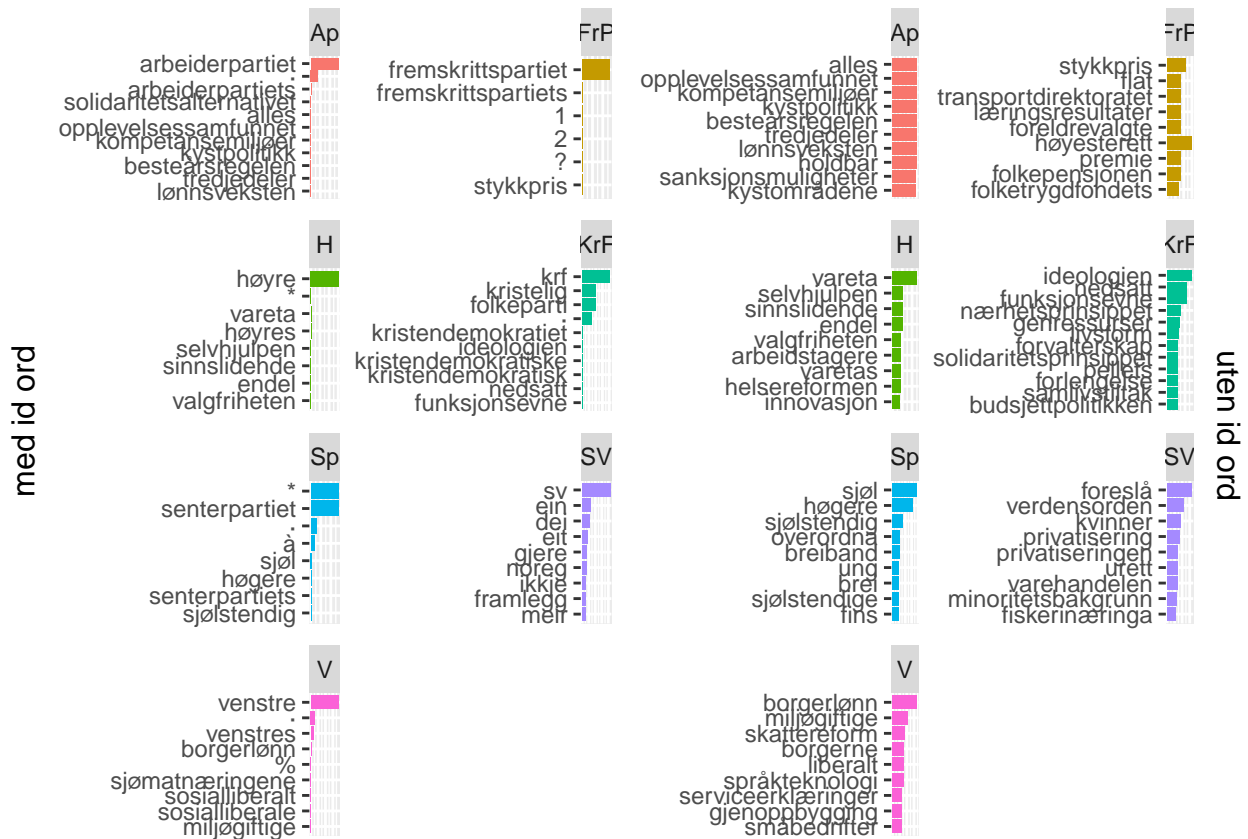
ord2 %>%
  ggplot(aes(term, tf_idf, fill = parti)) +
  geom_col(show.legend = FALSE) +
  facet_wrap(~parti, ncol =2, scale = "free") +
  theme(axis.title.x=element_blank(),
        axis.text.x=element_blank(),
        axis.ticks.x=element_blank(),
        axis.title.y = element_blank()) +
  coord_flip() -> fig_fixed

# 1 poeng
library(gridExtra)

##
## Attaching package: 'gridExtra'
## The following object is masked from 'package:dplyr':
##
##   combine

grid.arrange(fig_unfixed,fig_fixed, ncol =2, left = " med id ord", right = "uten id ord")

```



Ordlister

Lag 2 lister som med ord som etter din mening skiller godt høyre-siden fra venstresiden. Lag deretter 2 tilsvarende lister som skiller for og mot globalisering. Bruk mellom 10 og 20 ord i hver av listene. Lag en ordliste og bruk denne til å plassere partiene i forhold til hverandre og over tid. Vis dette i en tabell hvor partiene er sortert fra lav til høy verdi på høyre - venstre dimensjonen.

```
# 1 poeng
ordliste_venstre_global <- dictionary(list(venstre = c("lønn", "kapital", "lønn", "kvinner", "urett", "privat", "velferd", "felleskap", "likestilling"),
                                           global = c("internasjonal", "globalisering", "handel", "integrasjon")))

parti_ren2 %>%
  cast_dfm(document, term, n) -> parti_dfm
dfm_v_g <- dfm(parti_dfm, dictionary = ordliste_venstre_global)
dfm_v_g
```

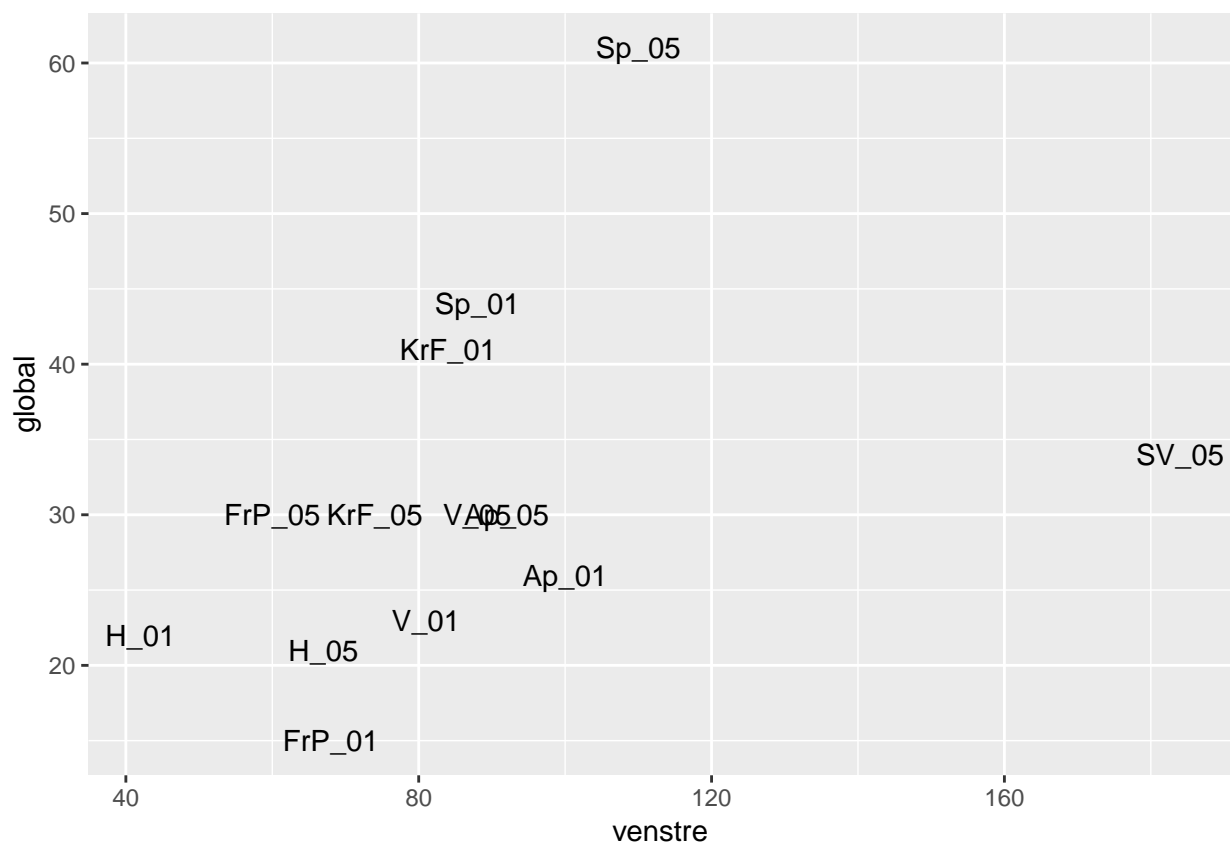
```
## Document-feature matrix of: 13 documents, 2 features (0.0% sparse).
## 13 x 2 sparse Matrix of class "dfm"
##           features
## docs  venstre global
## Ap_01    100    26
## V_01     81     23
## H_01     42     22
## KrF_01   84     41
## Sp_01   88     44
## FrP_01   68     15
## Ap_05   92     30
```

```
## SV_05      184    34
## V_05       88    30
## H_05       67    21
## KrF_05     74    30
## Sp_05     110    61
## FrP_05     60    30
```

Partiposisjoner figur

Lag en figur av dette i ggplot.

```
# 1 poeng
tmp <- convert(dfm_v_g, to = "data.frame")
ggplot(tmp, aes(x = venstre, y = global, label = document)) +
  geom_text()
```

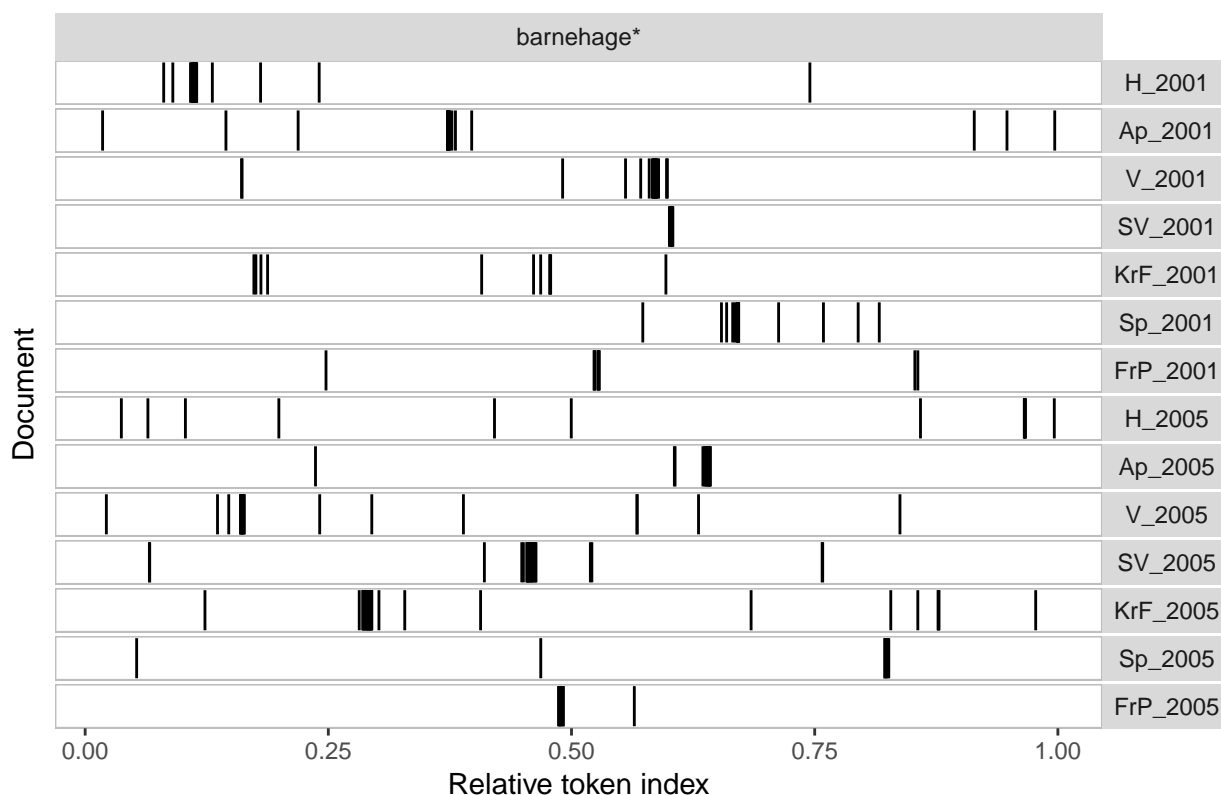


Bruk kwic til å danne deg en oversikt over parti-programmene

Nå skal du bruke kwic til å danne deg en oversikt over parti-programmene. Sammenlign hvordan partiene vektlegger barnehager. Vis med `textplot_xray` fra `quanteda`.

```
# 1 poeng
docnames(valgprogram) <- docvars(valgprogram, "navn")
kwic(valgprogram, pattern = phrase("barnehage*")) %>%
  textplot_xray()
```

Lexical dispersion plot



Cosine

Bruk cosine til å kalkulere forskjeller mellom tekstene. Vis dette som en tabell med 3 desimaler på hvert tall.

```
# 1 poeng
textstat_simil(parti_dfm, margin = "document", method = "cosine") %>% # regn ut cosine likhet
  round(3) #
```

```
##      Ap_01  V_01  H_01 KrF_01 Sp_01 FrP_01 Ap_05 SV_05  V_05  H_05
## V_01  0.991
## H_01  0.982 0.983
## KrF_01 0.987 0.987 0.980
## Sp_01  0.992 0.991 0.977  0.981
## FrP_01 0.974 0.977 0.978  0.965 0.976
## Ap_05  0.989 0.985 0.980  0.984 0.987  0.970
## SV_05  0.991 0.987 0.985  0.979 0.989  0.978 0.988
## V_05   0.989 0.994 0.991  0.987 0.988  0.983 0.986 0.991
## H_05   0.984 0.982 0.989  0.974 0.983  0.972 0.983 0.988 0.987
## KrF_05 0.989 0.989 0.977  0.991 0.990  0.971 0.989 0.985 0.987 0.980
## Sp_05  0.992 0.991 0.982  0.986 0.996  0.979 0.989 0.991 0.991 0.983
## FrP_05 0.981 0.983 0.983  0.980 0.983  0.994 0.979 0.981 0.987 0.979
##      KrF_05 Sp_05
## V_01
## H_01
## KrF_01
## Sp_01
```

```
## FrP_01
## Ap_05
## SV_05
## V_05
## H_05
## KrF_05
## Sp_05 0.992
## FrP_05 0.983 0.985
```

Wordscores

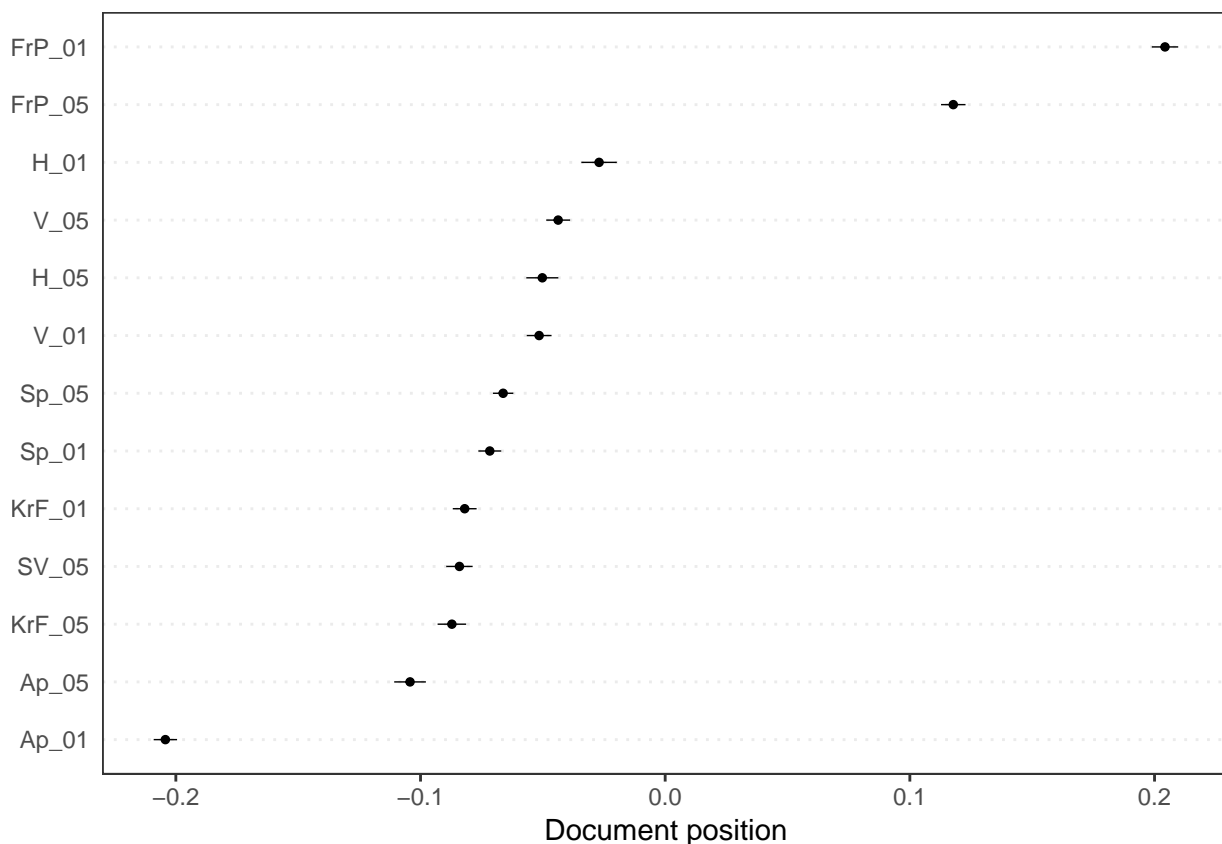
Ut fra hva du har sett til nå, og det du vet om norsk politikk, estimer en `wordscores` model som plasserer partiene langs høyre - venstre dimensjonen. Program for begge periodene skal være i samme figur.

```
# 1 poeng
```

```
ws <- textmodel_wordscores(x = parti_dfm, y = c(-1,NA, NA,NA,NA,1,NA,NA,NA,NA,NA,NA,NA)) # AP 01 = -1 o
parti_plassering <- predict(ws, se.fit = TRUE, level = .99, interval = "confidence", rescaling = "none")
```

```
## Warning: 16476 features in newdata not used in prediction.
```

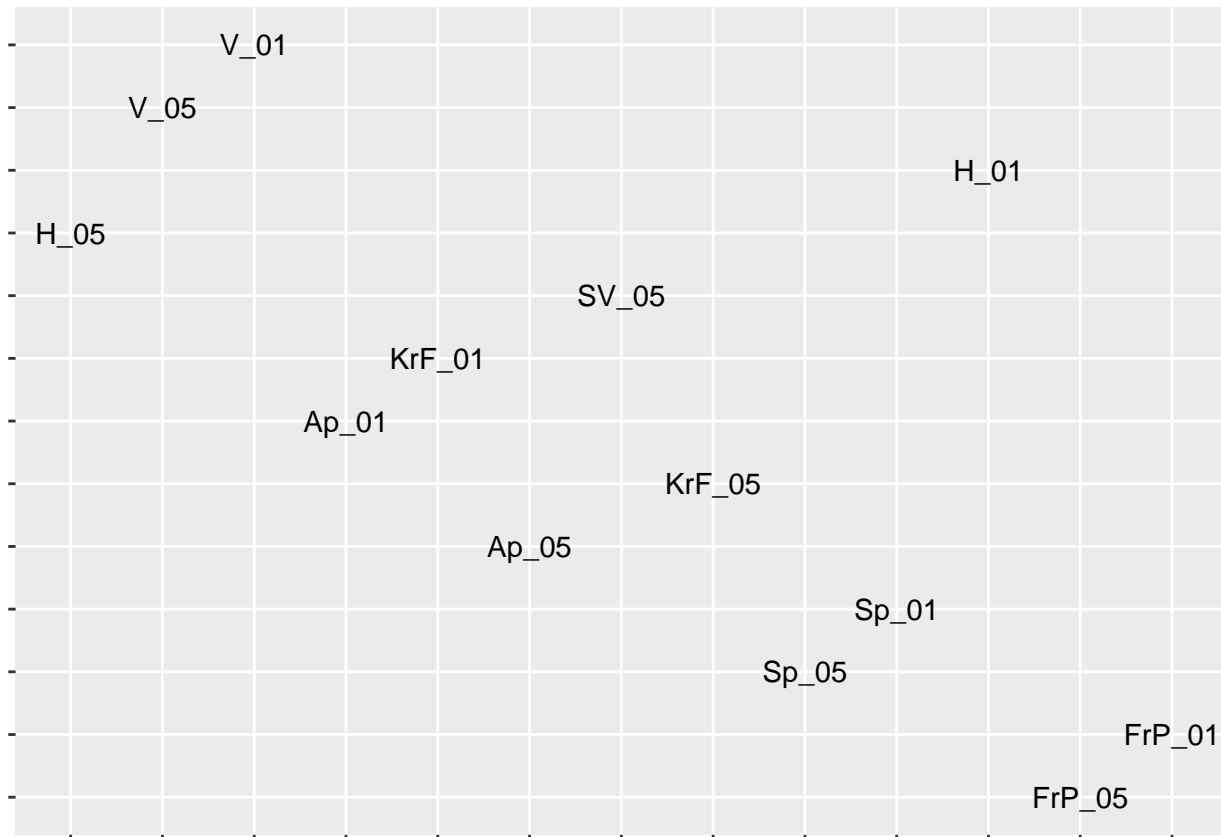
```
textplot_scale1d(parti_plassering, margin = "documents")
```



Korrespondanse analyse

Estimer en 2 dimensjonal korrespondanse analyse model og vis resultatene fra denne som en figur. Hva gir mest mening av denne modellen og wordscores?


```
# 1 poeng
ca2 <- textmodel_ca(parti_dfm, nd = 2)
ca2_pos <- as_tibble(cbind(ca2$rowcoord, names = ca2$rownames))
ggplot(ca2_pos, aes(x = Dim1, y = Dim2)) +
  geom_text(aes(label = names)) +
  theme(axis.title.x=element_blank(), # ta bort info å aksene
        axis.text.x=element_blank(),
        axis.title.y = element_blank(),
        axis.text.y=element_blank()) # Gir dette noe mening?
```



SKRIV MELLOM 100 og 200 ord om hvorfor du mener den ene modellen gir mer mening enn den andre her:

Spiller det noen rolle hvordan du preprosesserer tekstene?

Re-estimer den modellen du mener gir mest mening med de 2 spesifiseringene som `preText` viser at kan gi størst forskjell i resultater. Holder resultatene dine?

```
# 1 poeng
library(preText)

## preText: Diagnostics to Assess the Effects of Text Preprocessing Decisions
## Version 0.6.3 created on 2018-01-12.
## copyright (c) 2017, Matthew J. Denny, Penn State University
## Arthur Spirling, NYU
## Type vignette('getting_started_with_preText') to get started.
## Development website: https://github.com/matthewjdenny/preText
```

```
pp_corpus <- corpus(valgprogram)
pp_partier <- factorial_preprocessing(pp_corpus,
                                     infrequent_term_threshold = 0.12,
                                     parallel = TRUE, cores = 8) # dette tar noe tid
```

```
## Preprocessing 14 documents 128 different ways...
## Preprocessing documents 128 different ways on 8 cores. This may take a while...
## Preprocessing complete, loading in intermediate DFMs...
```

```
pp_partier$choices
```

##	removePunctuation	removeNumbers	lowercase	stem
## P-N-L-S-W-I-3	TRUE	TRUE	TRUE	TRUE
## N-L-S-W-I-3	FALSE	TRUE	TRUE	TRUE
## P-L-S-W-I-3	TRUE	FALSE	TRUE	TRUE
## L-S-W-I-3	FALSE	FALSE	TRUE	TRUE
## P-N-S-W-I-3	TRUE	TRUE	FALSE	TRUE
## N-S-W-I-3	FALSE	TRUE	FALSE	TRUE
## P-S-W-I-3	TRUE	FALSE	FALSE	TRUE
## S-W-I-3	FALSE	FALSE	FALSE	TRUE
## P-N-L-W-I-3	TRUE	TRUE	TRUE	FALSE
## N-L-W-I-3	FALSE	TRUE	TRUE	FALSE
## P-L-W-I-3	TRUE	FALSE	TRUE	FALSE
## L-W-I-3	FALSE	FALSE	TRUE	FALSE
## P-N-W-I-3	TRUE	TRUE	FALSE	FALSE
## N-W-I-3	FALSE	TRUE	FALSE	FALSE
## P-W-I-3	TRUE	FALSE	FALSE	FALSE
## W-I-3	FALSE	FALSE	FALSE	FALSE
## P-N-L-S-I-3	TRUE	TRUE	TRUE	TRUE
## N-L-S-I-3	FALSE	TRUE	TRUE	TRUE
## P-L-S-I-3	TRUE	FALSE	TRUE	TRUE
## L-S-I-3	FALSE	FALSE	TRUE	TRUE
## P-N-S-I-3	TRUE	TRUE	FALSE	TRUE
## N-S-I-3	FALSE	TRUE	FALSE	TRUE
## P-S-I-3	TRUE	FALSE	FALSE	TRUE
## S-I-3	FALSE	FALSE	FALSE	TRUE
## P-N-L-I-3	TRUE	TRUE	TRUE	FALSE
## N-L-I-3	FALSE	TRUE	TRUE	FALSE
## P-L-I-3	TRUE	FALSE	TRUE	FALSE
## L-I-3	FALSE	FALSE	TRUE	FALSE
## P-N-I-3	TRUE	TRUE	FALSE	FALSE
## N-I-3	FALSE	TRUE	FALSE	FALSE
## P-I-3	TRUE	FALSE	FALSE	FALSE
## I-3	FALSE	FALSE	FALSE	FALSE
## P-N-L-S-W-3	TRUE	TRUE	TRUE	TRUE
## N-L-S-W-3	FALSE	TRUE	TRUE	TRUE
## P-L-S-W-3	TRUE	FALSE	TRUE	TRUE
## L-S-W-3	FALSE	FALSE	TRUE	TRUE
## P-N-S-W-3	TRUE	TRUE	FALSE	TRUE
## N-S-W-3	FALSE	TRUE	FALSE	TRUE
## P-S-W-3	TRUE	FALSE	FALSE	TRUE
## S-W-3	FALSE	FALSE	FALSE	TRUE
## P-N-L-W-3	TRUE	TRUE	TRUE	FALSE
## N-L-W-3	FALSE	TRUE	TRUE	FALSE

## P-L-W-3	TRUE	FALSE	TRUE	FALSE
## L-W-3	FALSE	FALSE	TRUE	FALSE
## P-N-W-3	TRUE	TRUE	FALSE	FALSE
## N-W-3	FALSE	TRUE	FALSE	FALSE
## P-W-3	TRUE	FALSE	FALSE	FALSE
## W-3	FALSE	FALSE	FALSE	FALSE
## P-N-L-S-3	TRUE	TRUE	TRUE	TRUE
## N-L-S-3	FALSE	TRUE	TRUE	TRUE
## P-L-S-3	TRUE	FALSE	TRUE	TRUE
## L-S-3	FALSE	FALSE	TRUE	TRUE
## P-N-S-3	TRUE	TRUE	FALSE	TRUE
## N-S-3	FALSE	TRUE	FALSE	TRUE
## P-S-3	TRUE	FALSE	FALSE	TRUE
## S-3	FALSE	FALSE	FALSE	TRUE
## P-N-L-3	TRUE	TRUE	TRUE	FALSE
## N-L-3	FALSE	TRUE	TRUE	FALSE
## P-L-3	TRUE	FALSE	TRUE	FALSE
## L-3	FALSE	FALSE	TRUE	FALSE
## P-N-3	TRUE	TRUE	FALSE	FALSE
## N-3	FALSE	TRUE	FALSE	FALSE
## P-3	TRUE	FALSE	FALSE	FALSE
## 3	FALSE	FALSE	FALSE	FALSE
## P-N-L-S-W-I	TRUE	TRUE	TRUE	TRUE
## N-L-S-W-I	FALSE	TRUE	TRUE	TRUE
## P-L-S-W-I	TRUE	FALSE	TRUE	TRUE
## L-S-W-I	FALSE	FALSE	TRUE	TRUE
## P-N-S-W-I	TRUE	TRUE	FALSE	TRUE
## N-S-W-I	FALSE	TRUE	FALSE	TRUE
## P-S-W-I	TRUE	FALSE	FALSE	TRUE
## S-W-I	FALSE	FALSE	FALSE	TRUE
## P-N-L-W-I	TRUE	TRUE	TRUE	FALSE
## N-L-W-I	FALSE	TRUE	TRUE	FALSE
## P-L-W-I	TRUE	FALSE	TRUE	FALSE
## L-W-I	FALSE	FALSE	TRUE	FALSE
## P-N-W-I	TRUE	TRUE	FALSE	FALSE
## N-W-I	FALSE	TRUE	FALSE	FALSE
## P-W-I	TRUE	FALSE	FALSE	FALSE
## W-I	FALSE	FALSE	FALSE	FALSE
## P-N-L-S-I	TRUE	TRUE	TRUE	TRUE
## N-L-S-I	FALSE	TRUE	TRUE	TRUE
## P-L-S-I	TRUE	FALSE	TRUE	TRUE
## L-S-I	FALSE	FALSE	TRUE	TRUE
## P-N-S-I	TRUE	TRUE	FALSE	TRUE
## N-S-I	FALSE	TRUE	FALSE	TRUE
## P-S-I	TRUE	FALSE	FALSE	TRUE
## S-I	FALSE	FALSE	FALSE	TRUE
## P-N-L-I	TRUE	TRUE	TRUE	FALSE
## N-L-I	FALSE	TRUE	TRUE	FALSE
## P-L-I	TRUE	FALSE	TRUE	FALSE
## L-I	FALSE	FALSE	TRUE	FALSE
## P-N-I	TRUE	TRUE	FALSE	FALSE
## N-I	FALSE	TRUE	FALSE	FALSE
## P-I	TRUE	FALSE	FALSE	FALSE
## I	FALSE	FALSE	FALSE	FALSE

## P-N-L-S-W	TRUE	TRUE	TRUE	TRUE
## N-L-S-W	FALSE	TRUE	TRUE	TRUE
## P-L-S-W	TRUE	FALSE	TRUE	TRUE
## L-S-W	FALSE	FALSE	TRUE	TRUE
## P-N-S-W	TRUE	TRUE	FALSE	TRUE
## N-S-W	FALSE	TRUE	FALSE	TRUE
## P-S-W	TRUE	FALSE	FALSE	TRUE
## S-W	FALSE	FALSE	FALSE	TRUE
## P-N-L-W	TRUE	TRUE	TRUE	FALSE
## N-L-W	FALSE	TRUE	TRUE	FALSE
## P-L-W	TRUE	FALSE	TRUE	FALSE
## L-W	FALSE	FALSE	TRUE	FALSE
## P-N-W	TRUE	TRUE	FALSE	FALSE
## N-W	FALSE	TRUE	FALSE	FALSE
## P-W	TRUE	FALSE	FALSE	FALSE
## W	FALSE	FALSE	FALSE	FALSE
## P-N-L-S	TRUE	TRUE	TRUE	TRUE
## N-L-S	FALSE	TRUE	TRUE	TRUE
## P-L-S	TRUE	FALSE	TRUE	TRUE
## L-S	FALSE	FALSE	TRUE	TRUE
## P-N-S	TRUE	TRUE	FALSE	TRUE
## N-S	FALSE	TRUE	FALSE	TRUE
## P-S	TRUE	FALSE	FALSE	TRUE
## S	FALSE	FALSE	FALSE	TRUE
## P-N-L	TRUE	TRUE	TRUE	FALSE
## N-L	FALSE	TRUE	TRUE	FALSE
## P-L	TRUE	FALSE	TRUE	FALSE
## L	FALSE	FALSE	TRUE	FALSE
## P-N	TRUE	TRUE	FALSE	FALSE
## N	FALSE	TRUE	FALSE	FALSE
## P	TRUE	FALSE	FALSE	FALSE
##	FALSE	FALSE	FALSE	FALSE
##	removeStopwords	infrequent_terms	use_ngrams	
## P-N-L-S-W-I-3	TRUE	TRUE	TRUE	
## N-L-S-W-I-3	TRUE	TRUE	TRUE	
## P-L-S-W-I-3	TRUE	TRUE	TRUE	
## L-S-W-I-3	TRUE	TRUE	TRUE	
## P-N-S-W-I-3	TRUE	TRUE	TRUE	
## N-S-W-I-3	TRUE	TRUE	TRUE	
## P-S-W-I-3	TRUE	TRUE	TRUE	
## S-W-I-3	TRUE	TRUE	TRUE	
## P-N-L-W-I-3	TRUE	TRUE	TRUE	
## N-L-W-I-3	TRUE	TRUE	TRUE	
## P-L-W-I-3	TRUE	TRUE	TRUE	
## L-W-I-3	TRUE	TRUE	TRUE	
## P-N-W-I-3	TRUE	TRUE	TRUE	
## N-W-I-3	TRUE	TRUE	TRUE	
## P-W-I-3	TRUE	TRUE	TRUE	
## W-I-3	TRUE	TRUE	TRUE	
## P-N-L-S-I-3	FALSE	TRUE	TRUE	
## N-L-S-I-3	FALSE	TRUE	TRUE	
## P-L-S-I-3	FALSE	TRUE	TRUE	
## L-S-I-3	FALSE	TRUE	TRUE	
## P-N-S-I-3	FALSE	TRUE	TRUE	

## N-S-I-3	FALSE	TRUE	TRUE
## P-S-I-3	FALSE	TRUE	TRUE
## S-I-3	FALSE	TRUE	TRUE
## P-N-L-I-3	FALSE	TRUE	TRUE
## N-L-I-3	FALSE	TRUE	TRUE
## P-L-I-3	FALSE	TRUE	TRUE
## L-I-3	FALSE	TRUE	TRUE
## P-N-I-3	FALSE	TRUE	TRUE
## N-I-3	FALSE	TRUE	TRUE
## P-I-3	FALSE	TRUE	TRUE
## I-3	FALSE	TRUE	TRUE
## P-N-L-S-W-3	TRUE	FALSE	TRUE
## N-L-S-W-3	TRUE	FALSE	TRUE
## P-L-S-W-3	TRUE	FALSE	TRUE
## L-S-W-3	TRUE	FALSE	TRUE
## P-N-S-W-3	TRUE	FALSE	TRUE
## N-S-W-3	TRUE	FALSE	TRUE
## P-S-W-3	TRUE	FALSE	TRUE
## S-W-3	TRUE	FALSE	TRUE
## P-N-L-W-3	TRUE	FALSE	TRUE
## N-L-W-3	TRUE	FALSE	TRUE
## P-L-W-3	TRUE	FALSE	TRUE
## L-W-3	TRUE	FALSE	TRUE
## P-N-W-3	TRUE	FALSE	TRUE
## N-W-3	TRUE	FALSE	TRUE
## P-W-3	TRUE	FALSE	TRUE
## W-3	TRUE	FALSE	TRUE
## P-N-L-S-3	FALSE	FALSE	TRUE
## N-L-S-3	FALSE	FALSE	TRUE
## P-L-S-3	FALSE	FALSE	TRUE
## L-S-3	FALSE	FALSE	TRUE
## P-N-S-3	FALSE	FALSE	TRUE
## N-S-3	FALSE	FALSE	TRUE
## P-S-3	FALSE	FALSE	TRUE
## S-3	FALSE	FALSE	TRUE
## P-N-L-3	FALSE	FALSE	TRUE
## N-L-3	FALSE	FALSE	TRUE
## P-L-3	FALSE	FALSE	TRUE
## L-3	FALSE	FALSE	TRUE
## P-N-3	FALSE	FALSE	TRUE
## N-3	FALSE	FALSE	TRUE
## P-3	FALSE	FALSE	TRUE
## 3	FALSE	FALSE	TRUE
## P-N-L-S-W-I	TRUE	TRUE	FALSE
## N-L-S-W-I	TRUE	TRUE	FALSE
## P-L-S-W-I	TRUE	TRUE	FALSE
## L-S-W-I	TRUE	TRUE	FALSE
## P-N-S-W-I	TRUE	TRUE	FALSE
## N-S-W-I	TRUE	TRUE	FALSE
## P-S-W-I	TRUE	TRUE	FALSE
## S-W-I	TRUE	TRUE	FALSE
## P-N-L-W-I	TRUE	TRUE	FALSE
## N-L-W-I	TRUE	TRUE	FALSE
## P-L-W-I	TRUE	TRUE	FALSE

## L-W-I	TRUE	TRUE	FALSE
## P-N-W-I	TRUE	TRUE	FALSE
## N-W-I	TRUE	TRUE	FALSE
## P-W-I	TRUE	TRUE	FALSE
## W-I	TRUE	TRUE	FALSE
## P-N-L-S-I	FALSE	TRUE	FALSE
## N-L-S-I	FALSE	TRUE	FALSE
## P-L-S-I	FALSE	TRUE	FALSE
## L-S-I	FALSE	TRUE	FALSE
## P-N-S-I	FALSE	TRUE	FALSE
## N-S-I	FALSE	TRUE	FALSE
## P-S-I	FALSE	TRUE	FALSE
## S-I	FALSE	TRUE	FALSE
## P-N-L-I	FALSE	TRUE	FALSE
## N-L-I	FALSE	TRUE	FALSE
## P-L-I	FALSE	TRUE	FALSE
## L-I	FALSE	TRUE	FALSE
## P-N-I	FALSE	TRUE	FALSE
## N-I	FALSE	TRUE	FALSE
## P-I	FALSE	TRUE	FALSE
## I	FALSE	TRUE	FALSE
## P-N-L-S-W	TRUE	FALSE	FALSE
## N-L-S-W	TRUE	FALSE	FALSE
## P-L-S-W	TRUE	FALSE	FALSE
## L-S-W	TRUE	FALSE	FALSE
## P-N-S-W	TRUE	FALSE	FALSE
## N-S-W	TRUE	FALSE	FALSE
## P-S-W	TRUE	FALSE	FALSE
## S-W	TRUE	FALSE	FALSE
## P-N-L-W	TRUE	FALSE	FALSE
## N-L-W	TRUE	FALSE	FALSE
## P-L-W	TRUE	FALSE	FALSE
## L-W	TRUE	FALSE	FALSE
## P-N-W	TRUE	FALSE	FALSE
## N-W	TRUE	FALSE	FALSE
## P-W	TRUE	FALSE	FALSE
## W	TRUE	FALSE	FALSE
## P-N-L-S	FALSE	FALSE	FALSE
## N-L-S	FALSE	FALSE	FALSE
## P-L-S	FALSE	FALSE	FALSE
## L-S	FALSE	FALSE	FALSE
## P-N-S	FALSE	FALSE	FALSE
## N-S	FALSE	FALSE	FALSE
## P-S	FALSE	FALSE	FALSE
## S	FALSE	FALSE	FALSE
## P-N-L	FALSE	FALSE	FALSE
## N-L	FALSE	FALSE	FALSE
## P-L	FALSE	FALSE	FALSE
## L	FALSE	FALSE	FALSE
## P-N	FALSE	FALSE	FALSE
## N	FALSE	FALSE	FALSE
## P	FALSE	FALSE	FALSE
##	FALSE	FALSE	FALSE

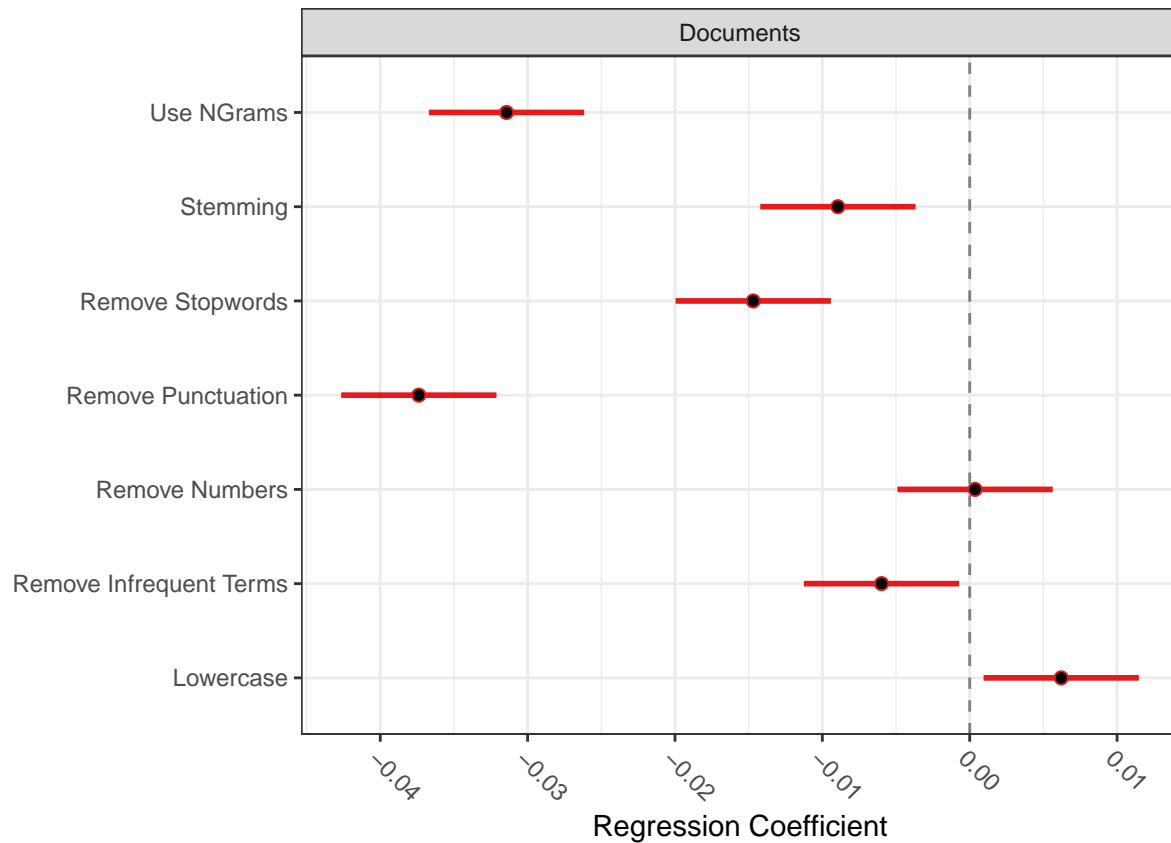
```

pt_resultat <- preText(
  pp_partier,
  num_comparisons = 15,
  verbose = FALSE,
  parallel = TRUE, cores = 8)

## Generating document distances...
## Generating preText Scores...
## Preprocessing documents 128 different ways on 8 cores. This may take a while...
## Paralellization complete! [1] 0.05863568 0.08726379 0.05834953 0.08649714 0.05862218 0.08525537
## [7] 0.05849530 0.07665209 0.05845751 0.09301101 0.05840892 0.09287874
## [13] 0.05841702 0.08723680 0.05858709 0.08704784 0.06654519 0.07751323
## [19] 0.06679084 0.07784527 0.05866267 0.07651981 0.05953731 0.07561009
## [25] 0.06553018 0.09547025 0.07006533 0.09564572 0.05892722 0.08631627
## [31] 0.06091945 0.08605982 0.07119372 0.08628118 0.06991956 0.08697495
## [37] 0.06962261 0.07903844 0.06404276 0.07137728 0.07906004 0.10415992
## [43] 0.07755372 0.10248353 0.07102095 0.09895260 0.07215743 0.09281935
## [49] 0.06936076 0.08973383 0.06978998 0.08849206 0.06880466 0.07918421
## [55] 0.06495789 0.07123151 0.08161106 0.10269949 0.07917072 0.10150092
## [61] 0.07746464 0.10275348 0.07851744 0.09321887 0.06423172 0.12141237
## [67] 0.06350556 0.12214124 0.06295756 0.11522244 0.06652359 0.11485531
## [73] 0.06306014 0.13737177 0.06333549 0.13790087 0.06304935 0.13071483
## [79] 0.06414534 0.12157434 0.10675143 0.13134651 0.10790411 0.13416748
## [85] 0.08940179 0.12732156 0.08187291 0.13694525 0.11792733 0.12261095
## [91] 0.12217633 0.12221142 0.10291275 0.14637188 0.10322049 0.16432081
## [97] 0.06374852 0.10898931 0.06321941 0.10919717 0.06206673 0.10738581
## [103] 0.06160512 0.11232048 0.06882626 0.14699007 0.06399687 0.14736260
## [109] 0.06137026 0.14439045 0.06094104 0.13776320 0.10055610 0.18382194
## [115] 0.09377497 0.19230105 0.07798294 0.14405032 0.07090757 0.14919825
## [121] 0.11141615 0.13958536 0.10847641 0.15222168 0.10705917 0.14013605
## [127] 0.10833603
## Generating regression results..
## The R^2 for this model is: 0.7591685
## Regression results (negative coefficients imply less risk):
##           Variable Coefficient    SE
## 1           Intercept         0.139 0.004
## 2 Remove Punctuation        -0.037 0.003
## 3 Remove Numbers             0.000 0.003
## 4 Lowercase                  0.006 0.003
## 5 Stemming                   -0.009 0.003
## 6 Remove Stopwords          -0.015 0.003
## 7 Remove Infrequent Terms   -0.006 0.003
## 8 Use NGrams                 -0.031 0.003
## Complete in: 57.904 seconds...

regression_coefficient_plot(pt_resultat, remove_intercept = TRUE)

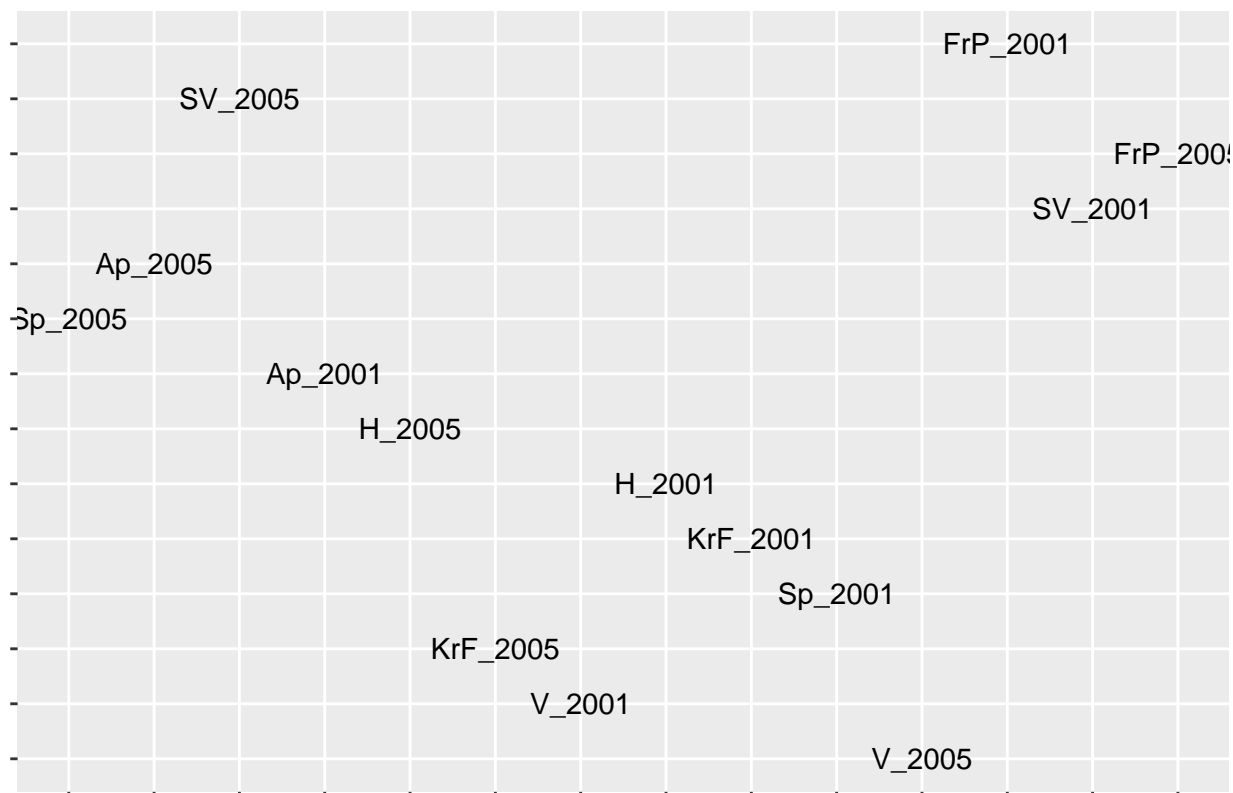
```



```
pp_corpus %>%
  tokens(what = "word",
         ngrams = 1:2, # inkluderer ngrams
         remove_punct = TRUE, # tar bort tegn
         remove_separators = TRUE) %>%
  dfm() -> ngram_dfm

ca2 <- textmodel_ca(ngram_dfm, nd = 2)
ca2_pos <- as_tibble(cbind(ca2$rowcoord, names = ca2$rownames))
ggplot(ca2_pos, aes(y = Dim1, x = Dim2)) + # snudd x og y skalane
  geom_text(aes(label = names)) +
  theme(axis.title.x=element_blank(), # ta bort info å aksene
        axis.text.x=element_blank(),
        axis.title.y = element_blank(),
        axis.text.y=element_blank()) + # Gir dette noe mening?
  ggtitle("Hva er den andre (første) dimensjonen?")
```


Hva er den andre (første) dimensjonen?



SKRIV EGEN VURDERING HER:

3 poeng

**TEST MED KNIT AT SKRIPTET KJØRER UTEN FEIL FØR DU LEVERER
INN**