# Topic 6: Topic Analysis Homework

Steven Cognac

2022-05-10

## Load the data

```
comments_df <- read_csv("https://raw.githubusercontent.com/MaRo406/EDS_231-text-sentiment/main/dat/comm
```

Now we'll build and clean the corpus

```
epa_corp <- corpus(x = comments_df, text_field = "text")
epa_corp.stats <- summary(epa_corp)
```

Now let's tokenize our dataset and remove stopwords

```
toks <- tokens(epa_corp, remove_punct = TRUE, remove_numbers = TRUE)
#I added some project-specific stop words here
add_stops <- c(stopwords("en"),"environmental", "justice", "ej", "epa", "public", "comment")

toks1 <- tokens_select(toks, pattern = add_stops, selection = "remove")
```

And now convert to a document-feature matrix

```
dfm_comm<- dfm(toks1, tolower = TRUE)
dfm <- dfm_wordstem(dfm_comm)
dfm <- dfm_trim(dfm, min_docfreq = 2) #remove terms only appearing in one doc (min_termfreq = 10)

# remove rows (docs) with all zeros
sel_idx <- slam::row_sums(dfm) > 0
dfm <- dfm[sel_idx, ]
```
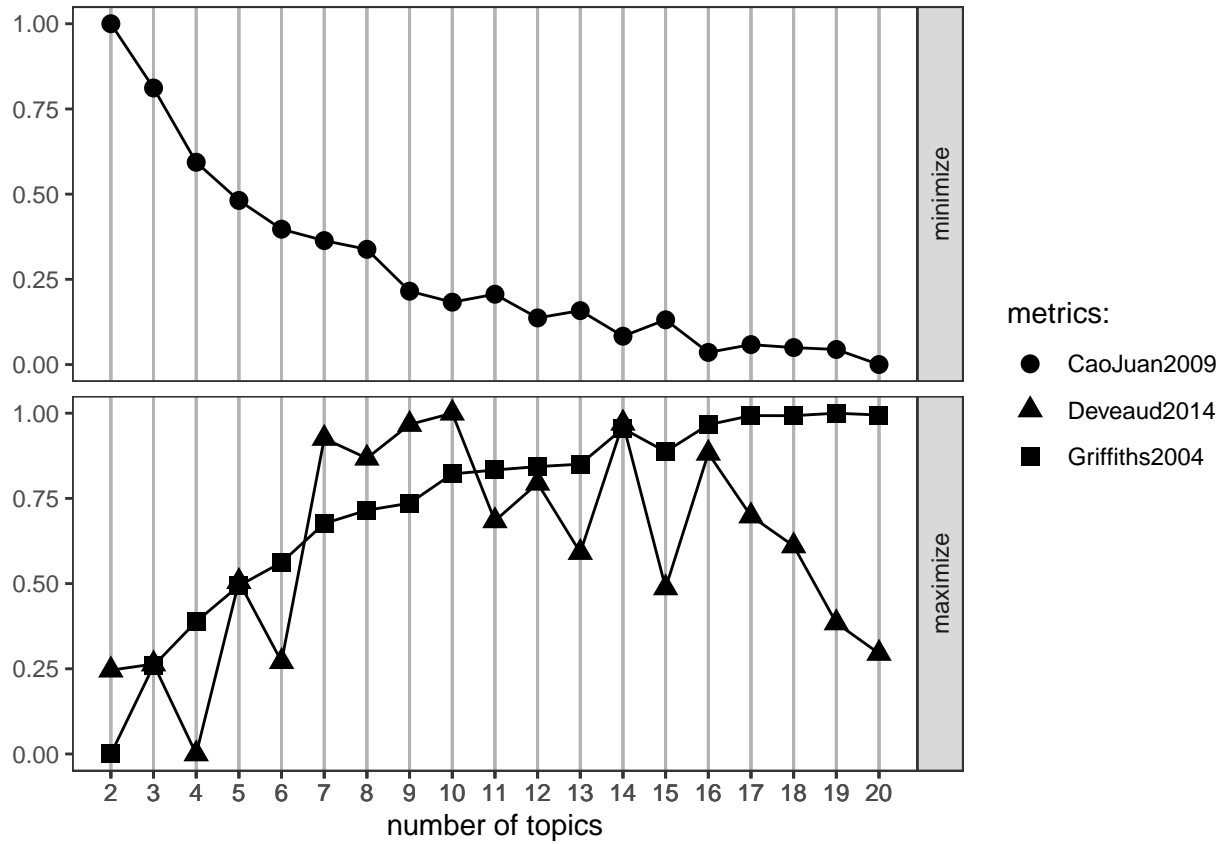
## Find optimal number of topics

```
#
result <- FindTopicsNumber(
  dfm,
  topics = seq(from = 2, to = 20, by = 1),
  metrics = c("CaoJuan2009",  "Deveaud2014", "Griffiths2004"),
  method = "Gibbs",
  control = list(seed = 77),
  verbose = TRUE
)
```

```
## fit models... done.
## calculate metrics:
##   CaoJuan2009... done.
##   Deveaud2014... done.
##   Griffiths2004... done.
```

```
FindTopicsNumber_plot(result)
```



# Selecting number of latent topics present from the comment letters

- Model 1 - **k=9** Latent Topics
- Model 2 - **k=4** Latent Topics
- Model 3 - **k=10** Latent topics
- Model 4 - **k=16** Latent topics

## Picking k values to check

Model 1 (k=9) is the original number of topics selected in the lab based on the 9 EPA priority areas: Rulemaking, Permitting, Compliance and Enforcement, Science, States and Local Governments, Federal Agencies, Community-based Work, Tribes and Indigenous People, National Measures. Model 2 (k=4) was chosen from reviewing the 2017-2020 EPA Environmental Justice Report's where I noticed there are 4 major goals/themes: (1) delivering environmental results; (2) cooperative federalism; (3) rule of law and fair process; and (4) building community capacity and engagement. Model's 3 (k=10) and 4 (k=16) were based

on the `FindTopicNumber()` optimization metrics. Based on the plots, 10 has the highest mximum value around where the minimization plot begins to descrease more slowely. In the minimization plot, 16 is where the plot begins to level out, indicating a stable number of topics is near. We also know 16 is a good option because there are 7 additional topics identified in the EPA's response to the public comments.

```r
# select number of topics
k9 <- 9
k4 <- 4
k10 <- 10
k16 <- 16

# run LDA model
topicModel_k9 <- LDA(dfm, k9, method="Gibbs", control=list(iter = 500, verbose = 25))
topicModel_k4 <- LDA(dfm, k4, method="Gibbs", control=list(iter = 500, verbose = 25))
topicModel_k10 <- LDA(dfm, k10, method="Gibbs", control=list(iter = 500, verbose = 25))
topicModel_k16 <- LDA(dfm, k16, method="Gibbs", control=list(iter = 500, verbose = 25))

#nTerms(dfm_comm)
tmResult9 <- posterior(topicModel_k9)
tmResult4 <- posterior(topicModel_k4)
tmResult10 <- posterior(topicModel_k10)
tmResult16 <- posterior(topicModel_k16)
# attributes(tmResult9)


# nTerms(dfm_comm)
# get beta from results
theta9 <- tmResult9$topics
theta4 <- tmResult4$topics
theta10 <- tmResult10$topics
theta16 <- tmResult16$topics

beta <- tmResult9$terms
```

Let's pull out the top 10 likelihood / probability of frequency in each topic for each latent topic number. Remember, just because we choose those number of topics, it doesn't mean LDS actually picked up on the correct topic.

```r
#terms(topicModel_k9, 10)
#terms(topicModel_k4, 10)
#terms(topicModel_k10, 10)
#terms(topicModel_k16, 10)
```

```r
# tidy terms
comment_topics9 <- tidy(topicModel_k9, matrix = "beta")
comment_topics4 <- tidy(topicModel_k4, matrix = "beta")
comment_topics10 <- tidy(topicModel_k10, matrix = "beta")
comment_topics16 <- tidy(topicModel_k16, matrix = "beta")

top_terms9 <- comment_topics9 %>%
  group_by(topic) %>%
  top_n(10, beta) %>%
  ungroup() %>%
```

```
  arrange(topic, -beta) %>%
  mutate(term = reorder(term, beta)) %>%
  ggplot(aes(term, beta, fill = factor(topic))) +
  geom_col(show.legend = FALSE) +
  facet_wrap(~ topic, scales = "free") +
  coord_flip() +
  labs(title = "k=9")

top_terms4 <- comment_topics4 %>%
  group_by(topic) %>%
  top_n(10, beta) %>%
  ungroup() %>%
  arrange(topic, -beta) %>%
  mutate(term = reorder(term, beta)) %>%
  ggplot(aes(term, beta, fill = factor(topic))) +
  geom_col(show.legend = FALSE) +
  facet_wrap(~ topic, scales = "free") +
  coord_flip() +
  labs(title = "k=4")

top_terms10 <- comment_topics10 %>%
  group_by(topic) %>%
  top_n(10, beta) %>%
  ungroup() %>%
  arrange(topic, -beta) %>%
  mutate(term = reorder(term, beta)) %>%
  ggplot(aes(term, beta, fill = factor(topic))) +
  geom_col(show.legend = FALSE) +
  facet_wrap(~ topic, scales = "free") +
  coord_flip() +
  labs(title = "k=10")

top_terms16 <- comment_topics16 %>%
  group_by(topic) %>%
  top_n(10, beta) %>%
  ungroup() %>%
  arrange(topic, -beta) %>%
  mutate(term = reorder(term, beta)) %>%
  ggplot(aes(term, beta, fill = factor(topic))) +
  geom_col(show.legend = FALSE) +
  facet_wrap(~ topic, scales = "free") +
  coord_flip() +
  labs(title = "k=16")

top_terms9
```
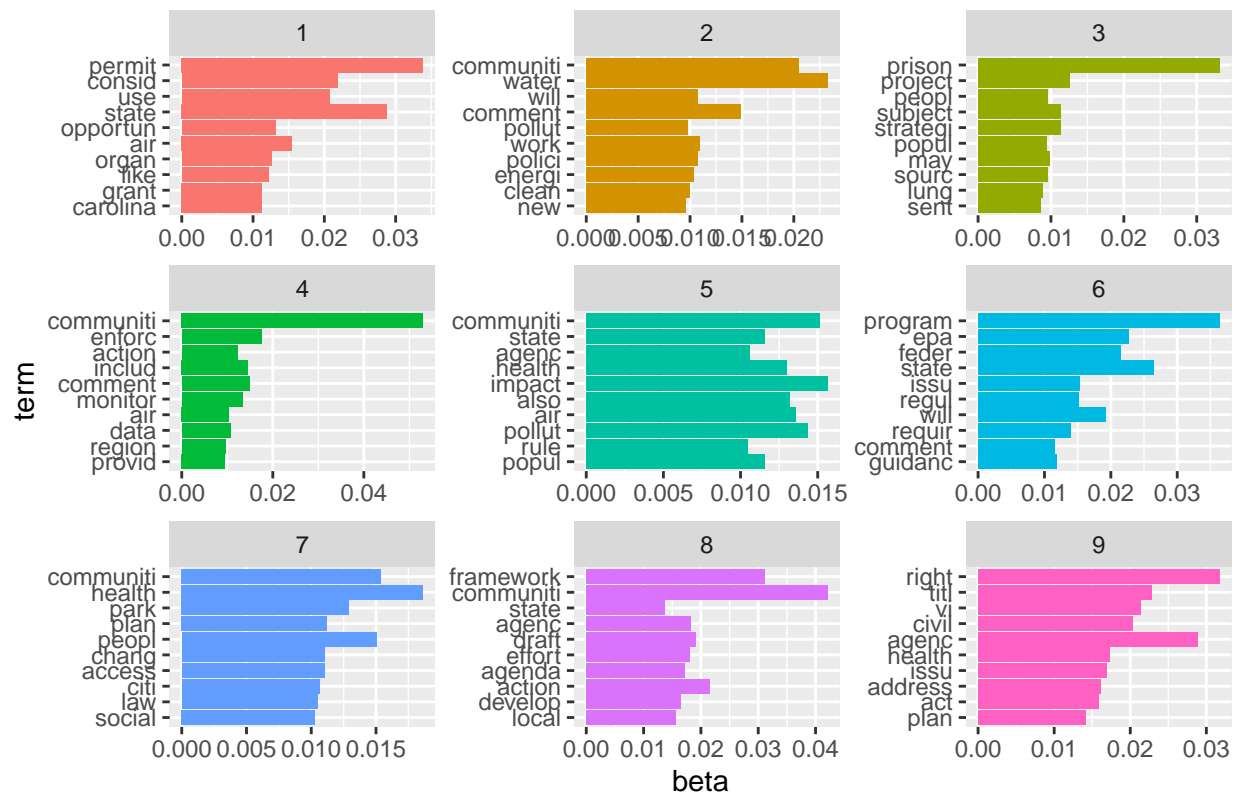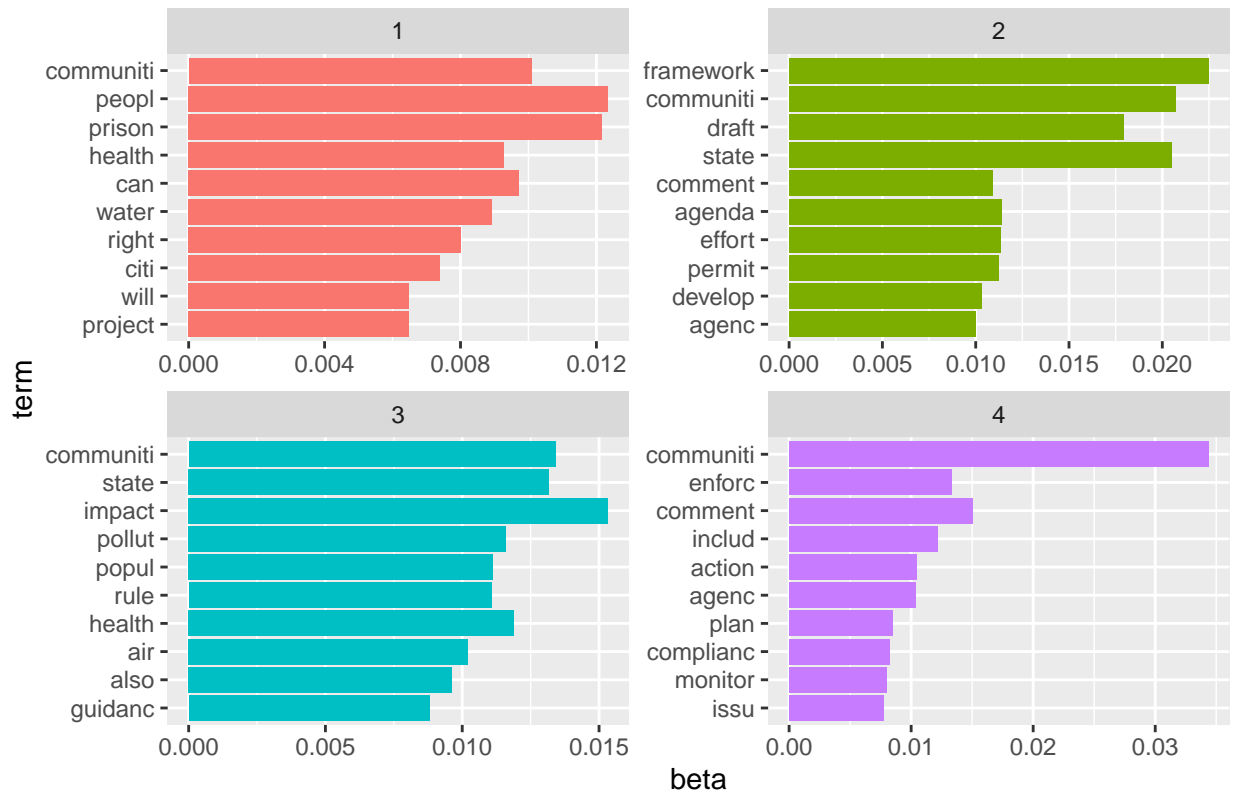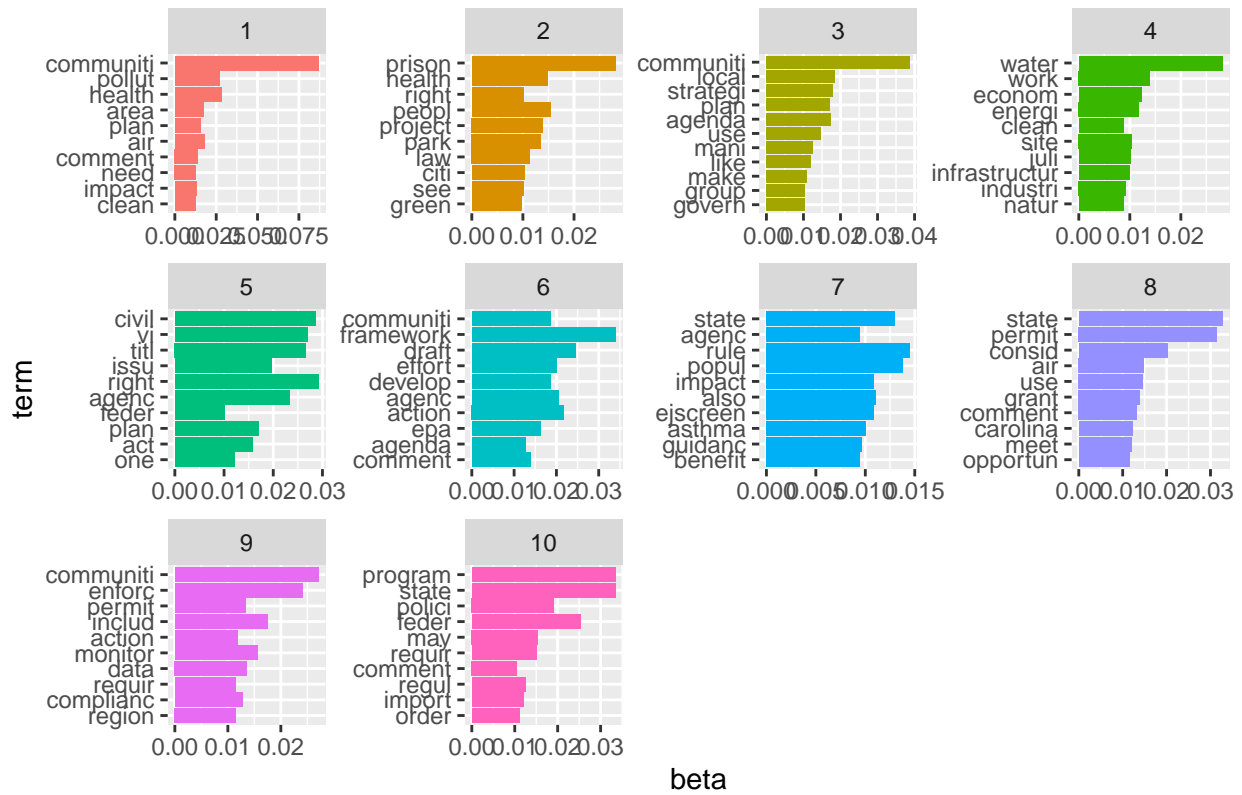
# k=9

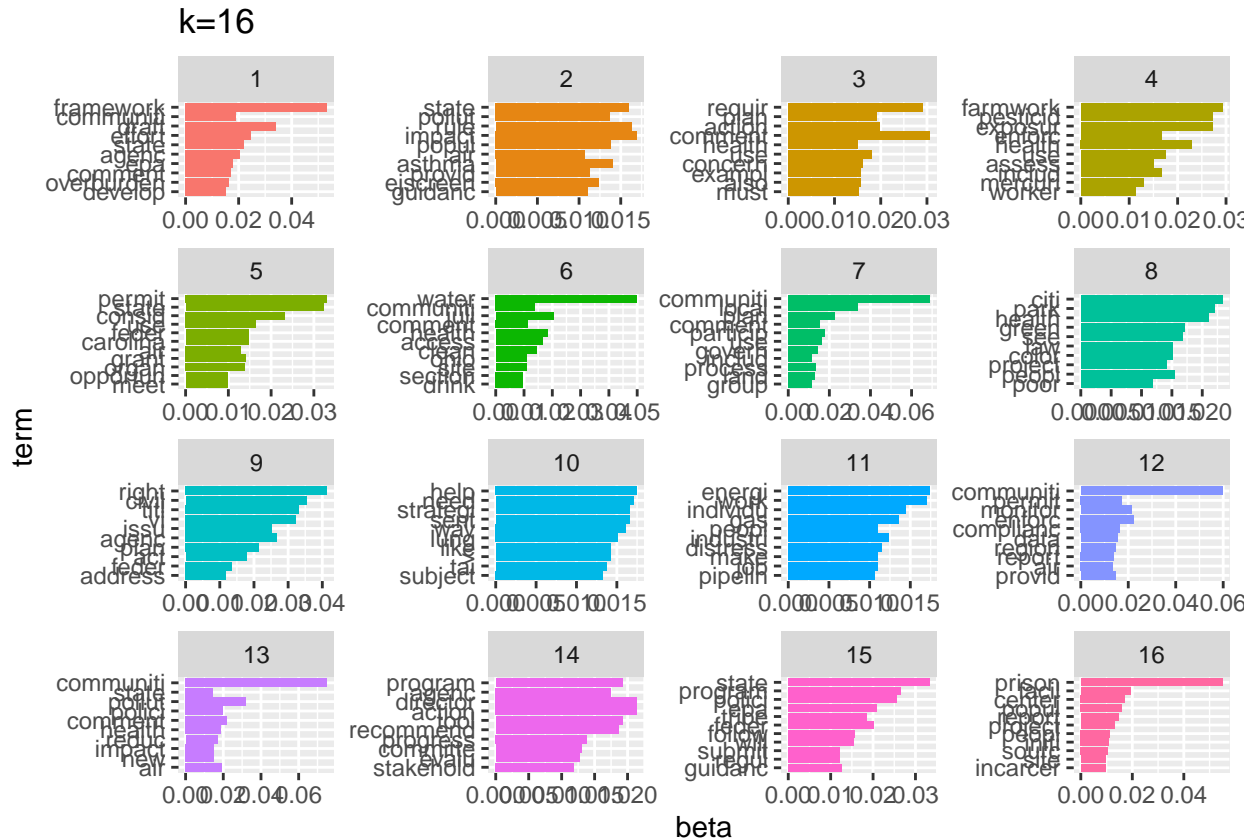

```
top_terms4
```

## k=4



```
top_terms10
```

# k=10



top_terms16

Let's assign names to the topics so we know what we are working with. We can name them by their top terms. Let's reassign the topic names to the top 5 words per topic.

```r
# re-assign names
top5termsPerTopic9 <- terms(topicModel_k9, 3)
top5termsPerTopic4 <- terms(topicModel_k4, 3)
top5termsPerTopic10 <- terms(topicModel_k10, 3)
top5termsPerTopic16 <- terms(topicModel_k16, 3)


# remove spaces
topicNames9 <- apply(top5termsPerTopic9, 2, paste, collapse=" ")
topicNames4 <- apply(top5termsPerTopic4, 2, paste, collapse=" ")
topicNames10 <- apply(top5termsPerTopic10, 2, paste, collapse=" ")
topicNames16 <- apply(top5termsPerTopic16, 2, paste, collapse=" ")
```

We can explore the theta matrix, which contains the distribution of each topic over each document

```r
exampleIds <- c(1, 2, 3, 4, 5, 6)
N <- length(exampleIds)

# get topic proportions from example documents
topicProportionExamples9 <- theta9[exampleIds,]
colnames(topicProportionExamples9) <- topicNames9
vizDataFrame9 <- melt(cbind(data.frame(topicProportionExamples9), document=factor(1:N)), variable.name =

topicProportionExamples4 <- theta4[exampleIds,]
```

```r
colnames(topicProportionExamples4) <- topicNames4
vizDataFrame4 <- melt(cbind(data.frame(topicProportionExamples4), document=factor(1:N)), variable.name =

topicProportionExamples10 <- theta10[exampleIds,]
colnames(topicProportionExamples10) <- topicNames10
vizDataFrame10 <- melt(cbind(data.frame(topicProportionExamples10), document=factor(1:N)), variable.name

topicProportionExamples16 <- theta16[exampleIds,]
colnames(topicProportionExamples16) <- topicNames16
vizDataFrame16 <- melt(cbind(data.frame(topicProportionExamples16), document=factor(1:N)), variable.name


# plot
terms9 <- ggplot(data = vizDataFrame9, aes(topic, value, fill = document), ylab = "proportion") +
  geom_bar(stat="identity") +
  theme(axis.text.x = element_text(angle = 90, hjust = 1)) +
  coord_flip() +
  facet_wrap(~ document, ncol = N) +
  labs(title = "K=9")

terms4 <- ggplot(data = vizDataFrame4, aes(topic, value, fill = document), ylab = "proportion") +
  geom_bar(stat="identity") +
  theme(axis.text.x = element_text(angle = 90, hjust = 1)) +
  coord_flip() +
  facet_wrap(~ document, ncol = N) +
  theme(legend.position="none") +
  labs(title = "k=4")

terms10 <- ggplot(data = vizDataFrame10, aes(topic, value, fill = document), ylab = "proportion") +
  geom_bar(stat="identity") +
  theme(axis.text.x = element_text(angle = 90, hjust = 1)) +
  coord_flip() +
  facet_wrap(~ document, ncol = N) +
  theme(legend.position="none") +
  labs(title = "k=10")

terms16 <- ggplot(data = vizDataFrame16, aes(topic, value, fill = document), ylab = "proportion") +
  geom_bar(stat="identity") +
  theme(axis.text.x = element_text(angle = 90, hjust = 1)) +
  coord_flip() +
  facet_wrap(~ document, ncol = N) +
  theme(legend.position="none") +
  labs(title = "k=16")

terms9
```
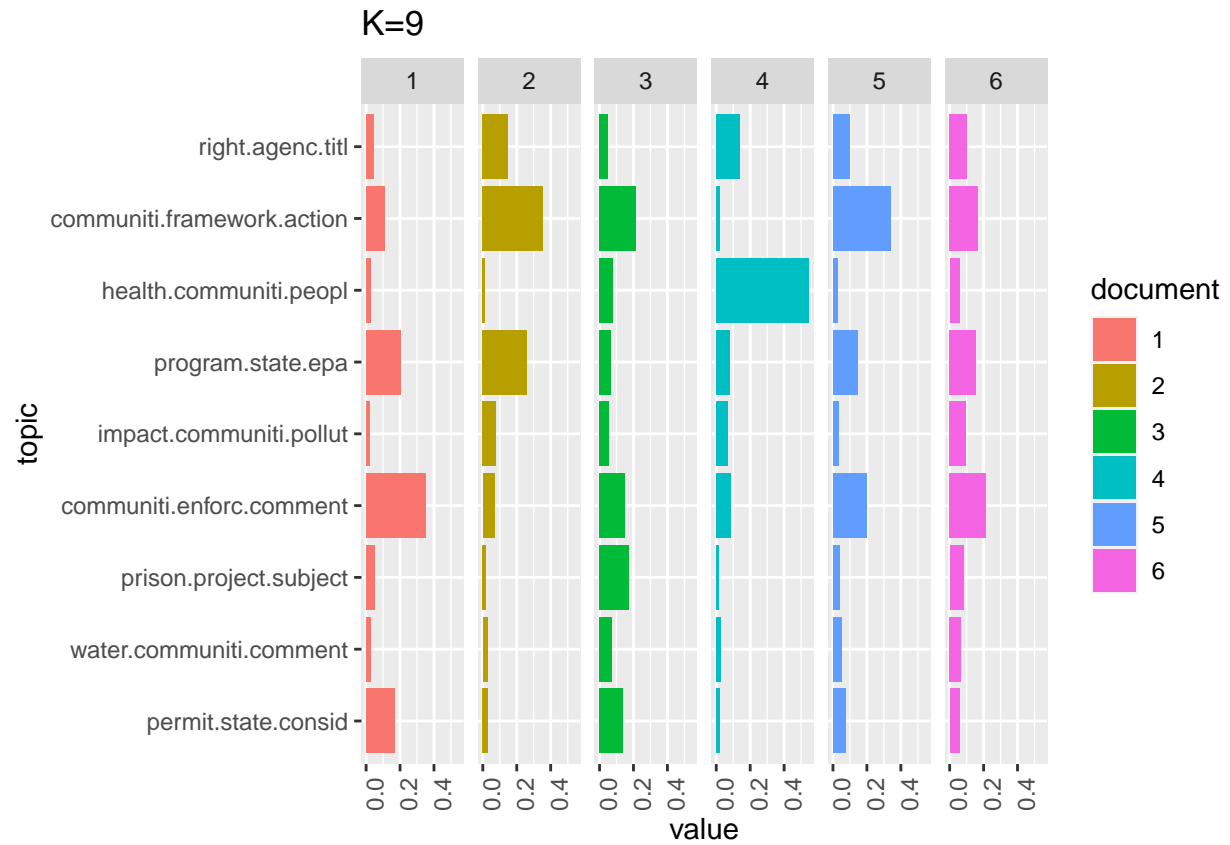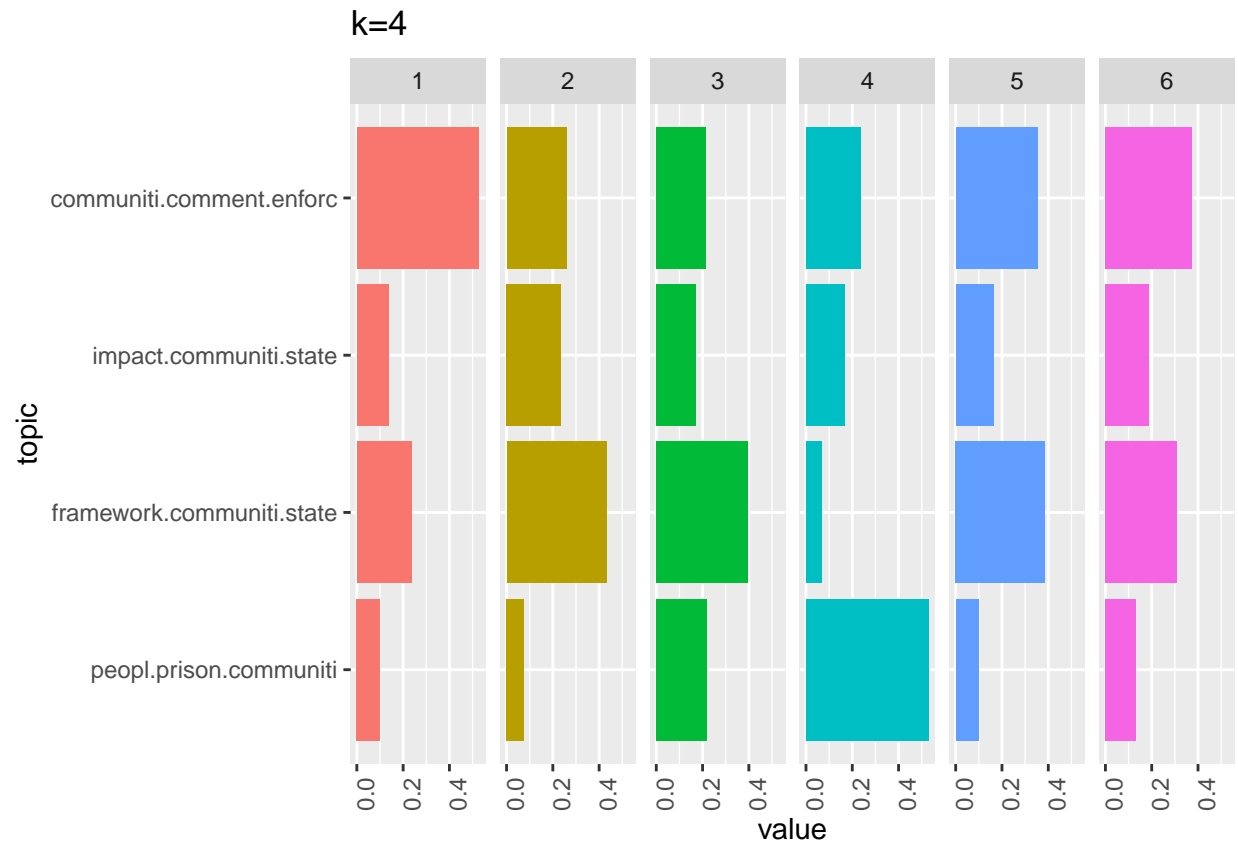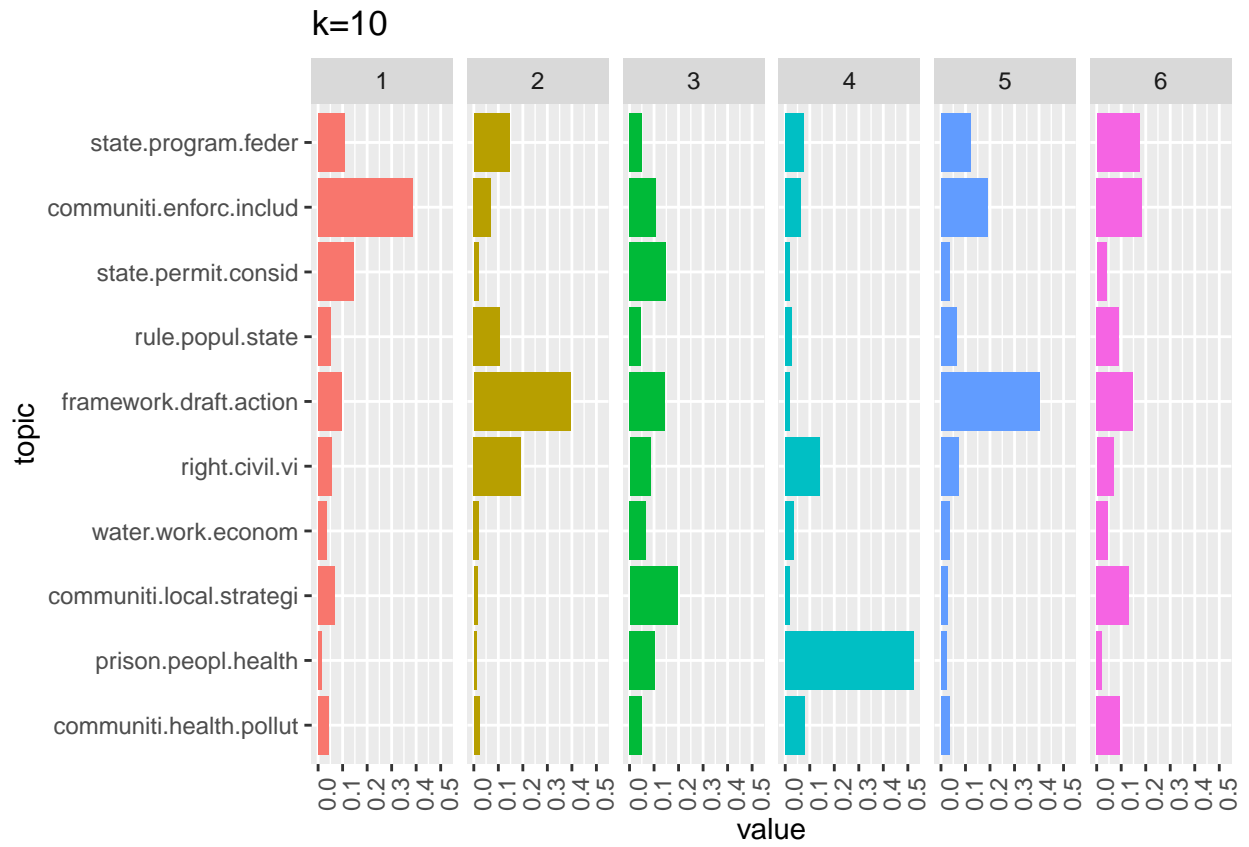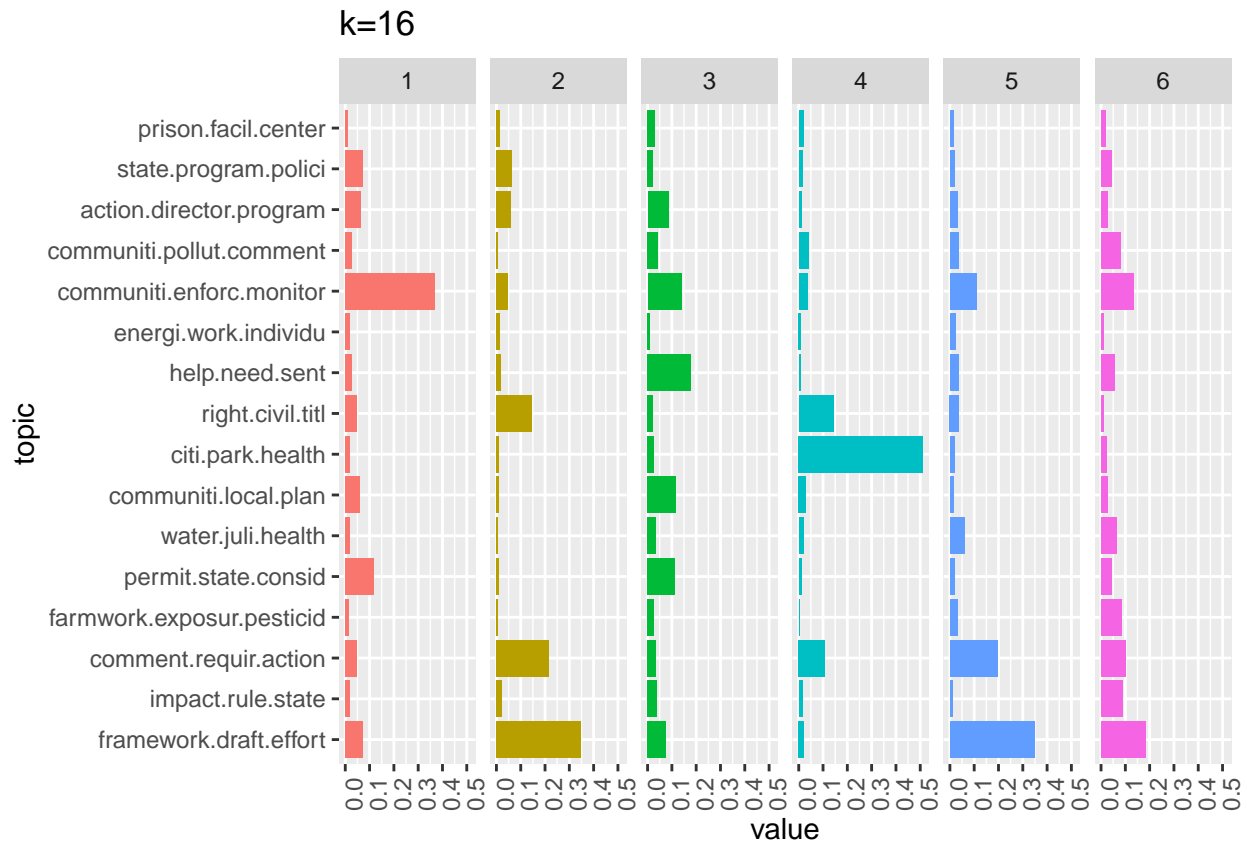
terms4

## k=4



terms10

k=10

terms16

k=16

```r
svd_tsne <- function(x){
  tsne(svd(x)$u)
}

json9 <- createJSON(
  phi = tmResult9$terms,
  theta = tmResult9$topics,
  doc.length = rowSums(dfm),
  vocab = colnames(dfm),
  term.frequency = colSums(dfm),
  mds.method = svd_tsne,
  plot.opts = list(xlab="",
                   ylab="")
)

json4 <- createJSON(
  phi = tmResult4$terms,
  theta = tmResult4$topics,
  doc.length = rowSums(dfm),
  vocab = colnames(dfm),
  term.frequency = colSums(dfm),
  mds.method = svd_tsne,
  plot.opts = list(xlab="",
                   ylab="")
)
```

```
json10 <- createJSON(
  phi = tmResult10$terms,
  theta = tmResult10$topics,
  doc.length = rowSums(dfm),
  vocab = colnames(dfm),
  term.frequency = colSums(dfm),
  mds.method = svd_tsne,
  plot.opts = list(xlab="",
                   ylab="")
)

json16 <- createJSON(
  phi = tmResult16$terms,
  theta = tmResult16$topics,
  doc.length = rowSums(dfm),
  vocab = colnames(dfm),
  term.frequency = colSums(dfm),
  mds.method = svd_tsne,
  plot.opts = list(xlab="",
                   ylab="")
)

#serVis(json9)
#serVis(json4)
#serVis(json10)
#serVis(json16)
```

## Picking the best model k-value

Based on the LDAvis and looking at the graphs of the top terms in each topic I would move forward with using 16 as the overall best value for k. When looking at the term plots for k=4, we get a lot of overlap between topics and minimual variation in unique topic subjects. When looking at k=10 as compared to k=16, there were again too much topic overlap with k=10. With the interactive LDAvis plots, the size of k=16 circles vary in size much more than k=10, indicating greater topic distriution.