

I love to develop and make things in R. Working on visualization styles, modeling techniques and general workflow problems.

(<https://www.hvitfeldt.me/blog/>)

(<https://www.hvitfeldt.me/GenArt/>)

(<https://www.hvitfeldt.me/about/>)

(<https://www.hvitfeldt.me/src/contrib/>)



/Emil_Hvitfeldt) /in/emilhvitfeldt/)
(https://github.com

/Email: Emil.Hvitfeldt@gmail.com | Template by <https://www.b5c.net/>

(<https://bootstrapious.com/free-templates>) & ported to Hugo by Kishan B (<https://github.com/kishaningithub>)

> Working showcase

- n) apply our model to the data to predict the author of the disputed papers

10/29/19, 9:12 PM

EMIL HVITFELDT
(HTTPS://WWW.HVITFELDT.ME/)

The Interesting thing in this was that the authorship of these papers were not consistent. In This is where we come in, in this blog post will we try to see if we are able to classify the troublesome papers.

- If you would like to read more about this story including past attempts to solve this problem please read How Statistics Solved a 175-Year-Old Mystery About Alexander Hamilton (https://priceconomics.com/how-statistics-solved-a-175-year-old-mystery-about/) by Ben Christopher (https://www.hvitfeldt.me/GenArt/)

- > About Libraries (https://www.hvitfeldt.me/about/)
- We will start by loading the libraries which includes
- > ERANmnet that will be used to construct the predictive model (https://www.hvitfeldt.me/src/contrib/)

```
library(tidyverse)
library(tidytext)
library(gutenbergr)
library(glmnet)
library(marginaleffects)
library(corrplot)
library(ggplot2)
```

/EmilHvitfeldt/in/emilhvitfeldt/ (https://github.com)

©2019 Emil Hvitfeldt
/EmilHvitfeldt | Template by Bootstrapious.com
We are lucky today because all of The Federalist Papers happens to be on gutenberg (https://bootstrapious.com/free-templates) & ported to Hugo by Kishan B (https://github.com/kishaningithub)

```
papers <- dutenberg_download(1404)
(HTTP://WWW.HVITFELDT.ME/)
head(papers, n = 10)
```

papers_sentences <- pull(papers, text)

str_c(collapse = " ") %>%

str_split(pattern = "\\.|\\?|\\!|") %

font_list() %>%

tibble(text = .) %>%

mutate(sentence = row_number())

We would like to assign each of these sentences to the corresponding article number and author. Thus we will first assign each of the 85 papers to the 3 authors and a group for the papers of interest.

EMIL HVITFELDT

(<https://www.hvitfeldt.me/>)

```
hamilton <- c(1, 6:9, 11:13, 15:17, 21)
madison <- c(10, 14, 18:20, 37:48)
jay <- c(2:5, 64)
unknown <- c(49:58, 62:63)
```

I love to develop and make things in R. Working on visualization styles, modeling techniques and general workflow problems.

Next we will simply look for lines that include "FEDERALIST No" as they would indicate the start of a paper and then label them accordingly.

> Blog

```
(https://www.hvitfeldt.me/papers_sentences %>%
  mutate(no = cumsum(str_detect(text,
```

> GenArt

```
(https://www.hvitfeldt.me/unnest_tokens(word, text) %>%
  mutate(author = case_when(no %in% ha
```

> About

```
(https://www.hvitfeldt.me/about/)
no %in% ma
no %in% ja
no %in% un
```

> ERAN

```
(https://www.hvitfeldt.me/erandata)
id = paste(no, sentence, sep
```

```
papers_words %>%
  count(author)
## # A tibble: 4 x 2
```

```
## author
## <chr> <int>
## 1 hamilton 114688
## 2 jay 8540
## 3 madison 45073
## 4 unknown 24471
```

```
/Emil Hvitfeldt | Template by
@EmilHvitfeldt | Bootstrapious.com
```

We see that Jay didn't post as many articles as the other two gentlemen so we will exclude him from further analysis

```
papers_words <- papers_words %>%
  filter(author != "jay")
```

Predictive modeling

I love to develop and make things in R. Working on visualization styles, modeling techniques and general workflow problems.

To make this predictive model we will use the term-

frequency matrix as our input and as the response will

be an indicator that Hamilton wrote the paper. For this

modeling we will use the `glmnet` package which fits a

generalized linear model via penalized maximum

likelihood. It is quite fast and works great with sparse

matrix input, hence the term-frequency matrix.

First we get the term-frequency matrix with the `cast_` family in `tidytext`.

```
papers_dtm <- papers_words %>%
  count(id, word, sort = TRUE) %>%
  cast_dtm(id, word, n)
```

We will need to define a response variable, which we

will do with `mutate`, along with an indicator

for our training set which will be the articles with

known authors.

```
meta <- data.frame(id = dimnames(paper
```

```
meta <- data.frame(id = dimnames(paper
  id, word, sort = TRUE) %>%
  cast_dtm(id, word, n)
```

```
meta <- data.frame(id = dimnames(paper
  id, word, sort = TRUE) %>%
  cast_dtm(id, word, n)
```

```
meta <- data.frame(id = dimnames(paper
  id, word, sort = TRUE) %>%
  cast_dtm(id, word, n)
```

We will use cross-validation to obtain the best value of the models tuning parameter. This part takes a couple of minutes.

EMIL HVITFELDT

(<https://www.hvitfeldt.me/>)

```
predictor <- papers_dtm[meta$train, ]
response <- meta$y[meta$train]
```

I love to develop and make

things in R. Working on

visualization styles, modeling

techniques and general

workflow problems.

After running the model, we will add the predicted values to our meta data.frame.

```
> Blog meta <- meta %>%
```

```
  mutate(pred = predict(model, newx =
    /blog/)
    s = model$lamb
```

```
> GenArt
```

It is now time to visualize the results. First we will look at how the training set have been separated.

```
/GenArt/)
```

```
> About
```

```
meta %>%
  (https://www.hvitfeldt.me
  filter(train) %>%
  /about/)
```

```
> ERAN
```

```
geom_boxplot(aes(fill = author)) +
  (https://www.hvitfeldt.me
  theme_minimal() +
  /src/contrib/
  labs(y = "predicted probability",
    x = "Article number") +
```

```
theme(legend.position = "top") +
```

```
scale_fill_manual(values = c("#30489
```

```
theme(axis.text.x = element_text(ang
```

```
(https://twitter.com/emilhvitfeldt)
(https://www.linkedin.com
```

/EmilHvitfeldt) /in/emilhv

(https://github.com

/EmilHvitfeldt). Template by

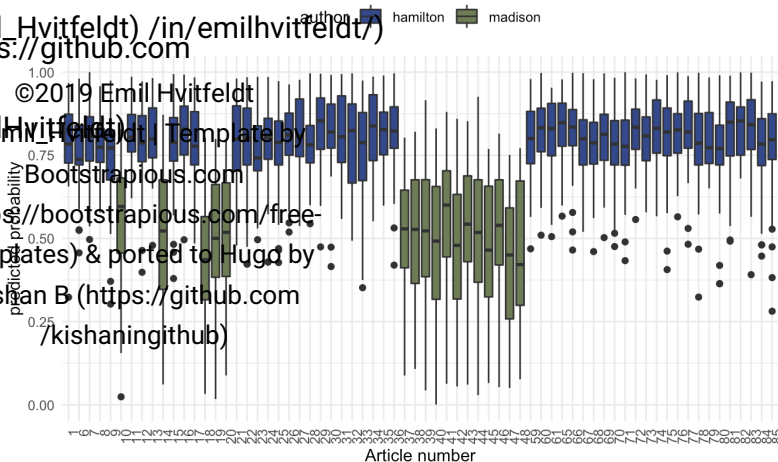
Bootstrapious.com

(https://bootstrapious.com/free-

templates) & ported to Hugo by

Kishan B (https://github.com

/kishaningithub)



The box plot of predicted probabilities, one value for

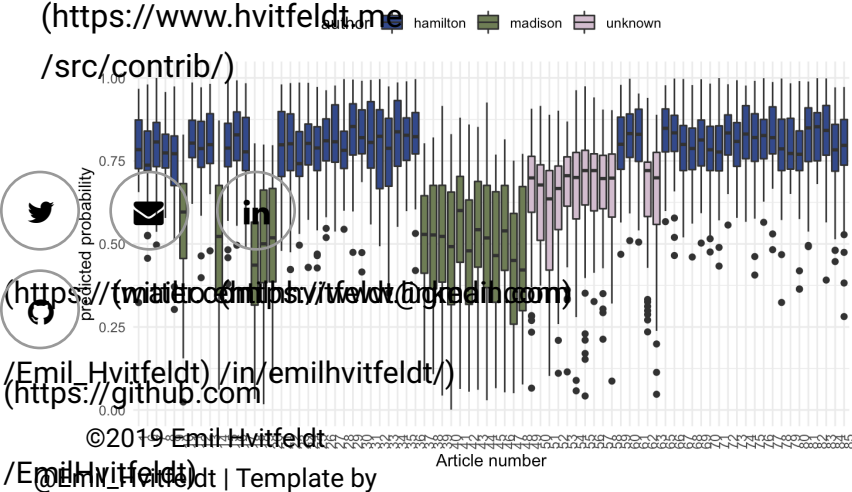
EMIL HVITFELDT

(<https://www.hvitfeldt.me/>)

each sentence for the 68 papers by Alexander Hamilton and James Madison. The probability represents the extent to which the model believe the sentence was written by Alexander Hamilton.

I love to develop and make things in R. Working on visualization styles, modeling Lets see if this model can settle the dispute of the 12 papers. We will plot the predicted probabilities of the unknown papers alongside the training set.

```
> Blog
(https://www.hvitfeldt.me
 /blog/) ggplot(aes(factor(no), pred)) +
  geom_boxplot(aes(fill = author)) +
  theme_minimal() +
  labs(y = "predicted probability",
       x = "Article number") +
  theme(legend.position = "top") +
  scale_fill_manual(values = c("#30489",
  /about/) theme(axis.text.x = element_text(ang
  > ERAN
(https://www.hvitfeldt.me
```



we notice that the predicted probabilities don't quite make up able to determine who the original author is. This can be due to a variety of different reasons. One of them could be that Madison wrote them and Hamilton edited them.

Despite the unsuccessful attempt to predict the secret author we still managed to showcase the method

EMIL HVITFELDT
(<https://www.hvitfeldt.me/>)

which while being unsuccessful in this case could provide useful in other cases.

I love to develop and make things with. Working on visualization styles, modeling

Working showcase

techniques and general workflow problems.

Since the method proved unsuccessful in determining

the secret author did I decide to add an example where

> Blog the authorship is known. We will use the same data from

(<https://www.hvitfeldt.me/blog/>) earlier, only look at known Hamilton and Madison papers, train on some of them and show that the

algorithm is able to detect the authorship of the other.

> GenArt

(<https://www.hvitfeldt.me/GenArt/>)
papers_dtm <- papers_words %>%

filter(author != "unknown") %>%

> About

(<https://www.hvitfeldt.me/about/>)
count(id, word, sort = TRUE) %>%
cast_dtm(id, word, n)

> ERAs Are we let the first 16 papers that they wrote be the

(<https://www.hvitfeldt.me/src/contrib/>) test set and the rest be

meta <- data.frame(id = dimnames(paper

left_join(papers_words[!duplicated(p

mutate(y = as.numeric(author == "ham

train = no > 20)

Warning: Column `id` joining factor

character vector

©2019 Emil Hvitfeldt

predictor <- papers_dtm[meta\$train,]

response <- meta\$y[meta\$train]

model <- cv.glmnet(predictor, response

meta <- meta %>%

mutate(pred = predict(model, newx =

s = model\$lamb

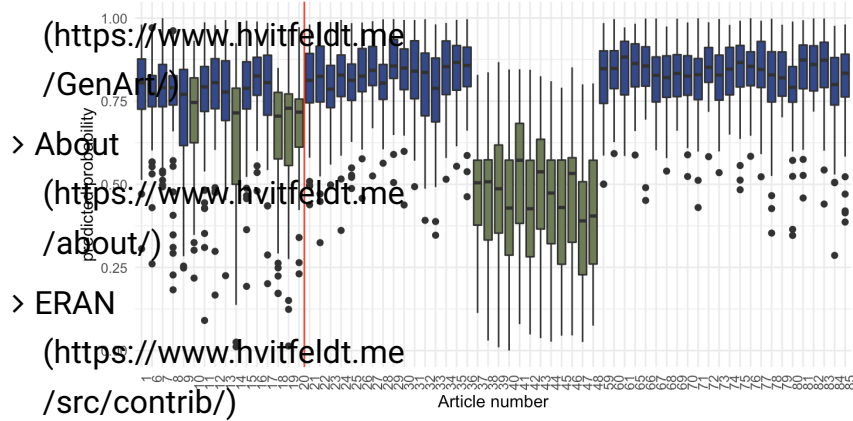


EMIL HVITFELDT

([HTTPS://WWW.HVITFELDT.ME/](https://www.hvitfeldt.me/))

```
ggplot(aes(factor(no), pred)) +  
  geom_boxplot(aes(fill = author)) +  
  theme_minimal() +  
  labs(y = "predicted probability",  
       x = "Article number") +  
  theme(legend.position = "top") +  
  scale_fill_manual(values = c("#30489  
> Blog theme(axis.text.x = element_text(ang  
(https://www.hvitfeldt.me/  
/blog/)
```

> GenArt



> About

(<https://www.hvitfeldt.me/about/>)

> ERAN

(<https://www.hvitfeldt.me/src/contrib/>)

So we see that while it isn't as crystal clear what what the test set predictions are giving us, they still give a pretty good indication.

([@emilhv/feldt](https://twitter.com/emilhv/feldt))

(<https://www.hvitfeldt.me/in/emilhv/feldt/>)

©2019 Emil Hvitfeldt

Emil Hvitfeldt | Template by @Emil_Hvitfeldt | Bootstrapious.com

(<https://bootstrapious.com/free-templates/>) & ported to Hugo by Kishan B (<https://github.com/kishaningithub>)

Sign in to comment

Styling with Markdown is supported

Sign in to comment

EMIL HVITFELDT
([HTTPS://WWW.HVITFELDT.ME/](https://www.hvitfeldt.me/))

I love to develop and make things in R. Working on visualization styles, modeling techniques and general workflow problems.

- > Blog
(<https://www.hvitfeldt.me/blog/>)
- > GenArt
(<https://www.hvitfeldt.me/GenArt/>)
- > About
(<https://www.hvitfeldt.me/about/>)
- > ERAN
(<https://www.hvitfeldt.me/src/contrib/>)



([@emilhvitfeldt](https://twitter.com/emilhvitfeldt)) (<mailto:emilhvitfeldt@gmail.com>) (<https://www.linkedin.com/company/emilhvitfeldt/>)
(<https://github.com/Emil-Hvitfeldt>)

©2019 Emil Hvitfeldt
/Emil-Hvitfeldt | Template by
Bootstrapious.com
(<https://bootstrapious.com/free-templates>) & ported to Hugo by
Kishan B (<https://github.com/kishaningithub>)