# Text Mining the History of Ideas

Mattias Skipper Rasmussen
filmsr@cas.au.dk

Frontiers in Digital Scholarship
Aarhus University
October 2016

## GENERIC "TOY" QUESTIONS

Some generic questions that may interest intellectual historians:

- Concept prevalence: how prevalent is concept $X$?
- Concept association: which concepts are typically associated with $X$?
- Topic discovery: which topics are frequently discussed?
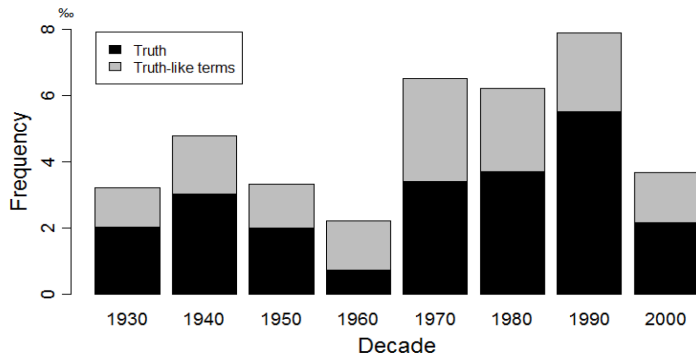- Topic prevalence: is there a correlation between the prevalence of different topics?

Intellectual historians typically assess such questions qualitatively.

I will illustrate how to approach them quantitatively.

## TEST CORPUS

- Sample of 324 philosophical articles from 1930-2010:
  - The Journal of Philosophy: 118 articles
  - The Philosophical Review: 206 articles

- Subcorpora:
  - 1930s: 38 articles
  - 1940s: 37 articles
    ⋮
  - 2010s: 53 articles

# PREVALENCE OF "TRUTH-LIKE" CONCEPTS (1)
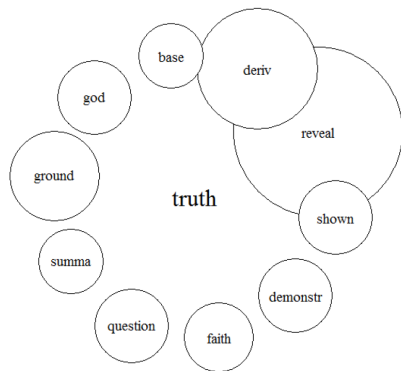
# PREVALENCE OF "TRUTH-LIKE" CONCEPTS (2)

```
146
147    ### Relative frequency of set of words over decades ###
148
149    # Function: words+decade --> rel. freq. of words in decade
150    words.freq.dec.f <- function(dtm,words.v,decade){
151        dec.art.m <- dtm[which(dtm[,1] %in% c(decade:(decade+9))),]
152        total <- sum(as.numeric(dec.art.m[,3:ncol(dec.art.m)]))
153        raw.words.freq <- sum(as.numeric(dec.art.m[,which(dimnames(dec.art.m)[[2]] %in% words.v)]))
154        words.freq.dec <- raw.words.freq/total
155        return(words.freq.dec)}
156
157    # Vector: rel. freq. of words over decades
158    words.freq.f <- function(dtm,words.v){
159        words.freq.l <- list()
160        for (i in 1:length(decades)){
161            words.freq.l <- c(words.freq.l,words.freq.dec.f(dtm,words.v,decades[[i]]))}
162        words.freq.v <- unlist(words.freq.l)
163        return(words.freq.v)}
164
165    # Plot: rel. freq. of words against decades
166    truth.v <- c("truth")
167    truth.phrases.v <- c("false","falsity","falsehood","truths","untrue","truthful")
168    barplot(words.freq.f(cor.sparse.md.m,truth.phrases.v),names.arg = decades,
169            main = "Frequency of 'truth/falsity' terms",xlab = "Decade",ylab = "Frequency")
170    barplot(words.freq.f(cor.sparse.md.m,truth.v),names.arg = decades,
171            main = "Frequency of 'truth/falsity' terms",xlab = "Decade",ylab = "Frequency",
172            legend("topright",legend = c("a", "b"),fill = c("black","grey")))
173
```
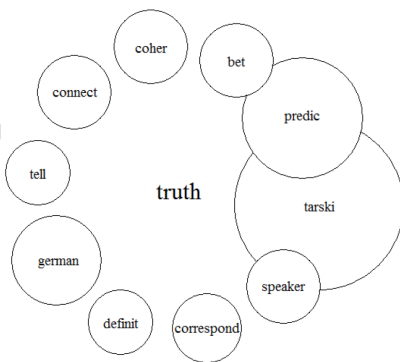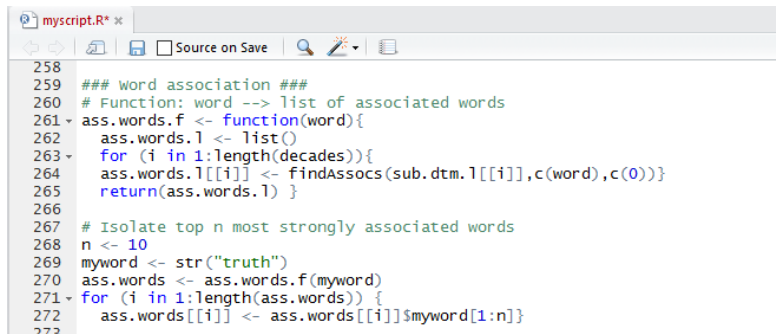
# WORDS ASSOCIATED WITH "TRUTH" (1)

# WORDS ASSOCIATED WITH "TRUTH" (2)

```
myscript.R* ×

      Source on Save

258
259   ### Word association ###
260   # Function: word --> list of associated words
261 ▾ ass.words.f <- function(word){
262     ass.words.l <- list()
263 ▾   for (i in 1:length(decades)){
264     ass.words.l[[i]] <- findAssocs(sub.dtm.l[[i]],c(word),c(0))}
265     return(ass.words.l) }
266
267   # Isolate top n most strongly associated words
268   n <- 10
269   myword <- str("truth")
270   ass.words <- ass.words.f(myword)
271 ▾ for (i in 1:length(ass.words)) {
272     ass.words[[i]] <- ass.words[[i]]$myword[1:n]}
273
```

# TOPIC DISCOVERY (1)

Topic models are algorithms that allow us to identify hidden thematic patterns in (sets of) texts.
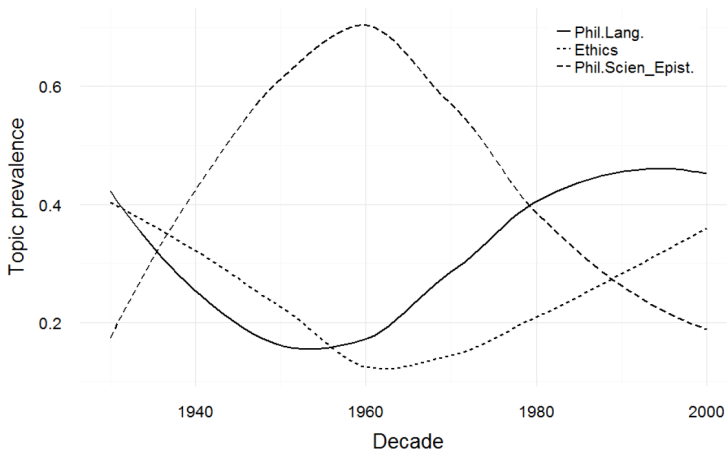
- A topic is treated as a distribution over words.
- A document is treated as a distribution over topics.

# TOPIC DISCOVERY (2)

myscript.R* ×

```
310
311   ### Topic prevalence ###
312
313   # Number of topics
314   n <- 3
315
316   # Run LDA using Gibbs sampling
317   ldaOut <- LDA(cor.sparse.dtm,n,method="Gibbs",control = list(seed=seed))
318
319   # Posterior probability of each word in each topic
320   ldaOutpost.1 <- posterior(ldaOut, cor.sparse.dtm)
321
322   # List of topics
323   top.wc.1 <- list()
324   for (i in 1:n){
325     top.wc.1[[i]] <- sort(ldaOutpost.1$terms[i,], decreasing = T)
326   }
327
328   # Word cloud of topic
329   top.nr <- 2
330   word.nr <- 25
331   greyscale <- brewer.pal(8,"Greys")
332   wordcloud(names(top.wc.1[[top.nr]][1:word.nr]),top.wc.1[[top.nr]][1:word.nr],
333             scale=c(8,.2), random.order=FALSE, rot.per=.15,colors = greyscale)
334
```

# TOPIC PREVALENCE

Thanks for your attention!