# News mining for conflict forecasting

Thomas Chadefaux — Uppsala

21 September 2017

# Sources

# Factiva searches

Can use booleans:

- Simple searches:
    - X and Y
    - X or Y
    - X not Y
    - More complex booleans:
        - X and (Y or Z)

- More complex:
    - "atleast2 deaths" $\rightarrow$ at least two mentions of "deaths"
    - "x same y" $\rightarrow$ looks for mentions of x in the same paragraph as y
    - "President w/3 Trump" $\rightarrow$ "president" within three words of Trump. E.g., President Donald J. Trump
    - "x/F100" $\rightarrow$ looks for x within first 100 words of the article
    - "fight*" $\rightarrow$ returns fighting, fighter, fights, etc.
    - Only look for keywords in articles of a certain size -E.g., "rebels and wc > 1000"

# Alternatives to Factiva

LexisNexis Newsdesk



Figure 1:
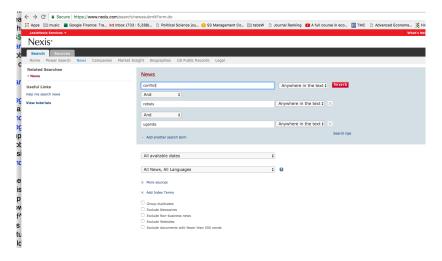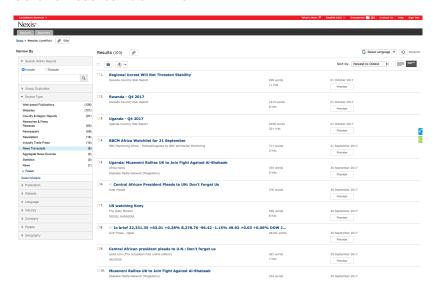
# What the results look like



Figure 2:

# Processing the results

# R package for factiva/lexisNexis

```r
# Install relevant package (a similar one exists for Facti
#install.packages('tm.plugin.lexisnexis')

library(tm) # text mining package
library(tm.plugin.lexisnexis)

# Import corpus
source <- LexisNexisSource("All_News,_All_Languages2017-09-
corpus <- Corpus(source, readerControl = list(language = N
```

# R package for factiva/lexisNexis

```r
# See how many articles were imported
corpus
```

```
## <<VCorpus>>
## Metadata:  corpus specific: 0, document level (indexed):
## Content:  documents: 55
```

```r
# Get an overview of the contents of the first article and
print(corpus[1]$content)
```

```
## [[1]]
## <<PlainTextDocument>>
## Metadata:  17
## Content:  chars: 2720
```

# R package for factiva/lexisNexis

```
# See the metadata for the first article
meta(corpus[[1]])
```

```
##   author      : character(0)
##   datetimestamp: 2017-09-20
##   description : character(0)
##   heading     : US watching Kony
##   id          : TheDaily201709201
##   language    : NA
##   origin      : The Daily Monitor
##   intro       : character(0)
##   section     : character(0)
##   subject     : character(0)
##   coverage    : character(0)
##   company     : character(0)
##   stocksymbol : character(0)
##   industry    : character(0)
##   type        : character(0)
```

# R package for factiva/lexisNexis

```r
# Get the text of the first article
corpus[[1]]$content
```

```
##  [1] "The US is working with African Union Mission in So
##  [2] "STUTTGART- The American troops have withdrawn from
##  [3] " leader Joseph Kony but they are \"keenly watching
##  [4] " are operating. The commander of the US Africa Com
## NA
##  [6] "The American government recently announced they we
##  [7] " outfit that has been weakened. Money spent Gen Wa
##  [8] " who had abducted and killed thousands in the regi
##  [9] " in 2006, fled to DR Congo and later to CAR in 200
## [10] " are now weak. \"This past April, it got to a poi
```

# R package for factiva/lexisNexis

```r
#remove stopwords using the standard list in tm
dtm <- DocumentTermMatrix(corpus)
inspect(dtm)
```

```
## <<DocumentTermMatrix (documents: 55, terms: 8046)>>
## Non-/sparse entries: 21138/421392
## Sparsity           : 95%
## Maximal term length: 55
## Weighting          : term frequency (tf)
## Sample             :
##                     Terms
## Docs                 and for from have said south that
##    DailyMoni201709109  56  11   12   13    2    17   14
##    DailyMoni201709125  24  10    6   12    0     0   12
##    TheDaily201709108   58  11   12   13    2    17   14
##    TheGuardi2017090219 100  34   24   14    2     0   51
##    TheIndepe201709163  34  12    7   14    7     1    8
##    TheNewYo2017082353  40   8    5    3    6     1   18
```

# Collect What?

# What to collect?

- Number of articles mentioning any keyword
- Number of hits for each keyword
    - More flexible
    - Multiple TS
- Store all possible articles about country x. Ideal for future uses

# Keywords

|  | Static | Dynamic |
|---|---|---|
| Human-defined | Chadefaux 2014 | |
| Automatically inferred | | APSR (forth.) |

Underlying question: build a data also useful for others, or maximize prediction power?

# More advanced processing?

- Tone
- Collocations/n-grams
- TABARI, etc.

Underlying questions: looking for perceptions, or event coding?

# Things to keep in mind

- Article inflation
- Articles correlation

# Vague paper idea

- Dynamic time warping for news count. Works well with irregular data.
  - Find /match patterns in TS data