

1 Identify languages by text analysis (Projekt 1)

1.1 A program that distinguishes Swedish from Norwegian

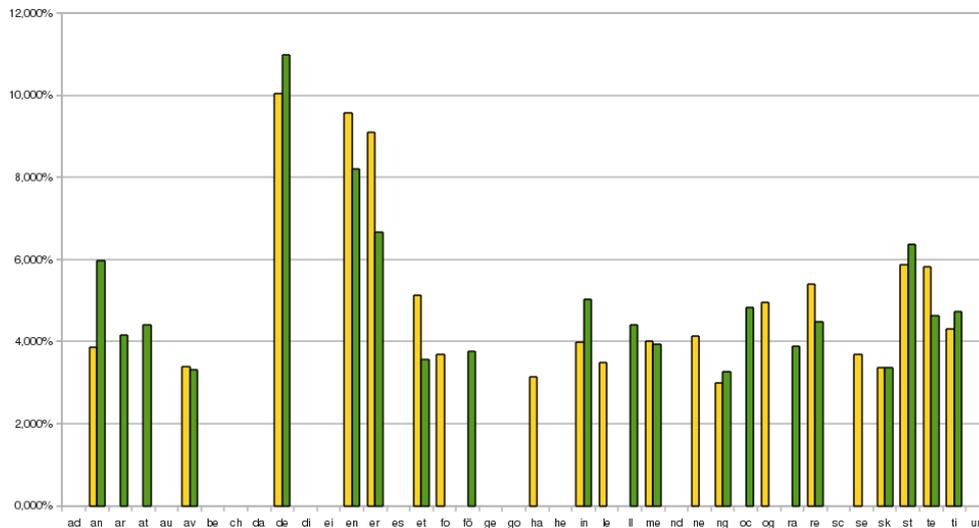
example texts As example texts to work with I have fetched ~20MB of data from Wikipedia using a shellscript that downloads the list of the “excellent articles” in the Norwegian and Swedish Wikipedia (currently ~200 articles each) and strips the relevant text out of the downloaded HTML-sourcecode (total word count of all wikipedia samples after processing: 2 million words)

In order to try the “extra challenge” I additionally fetched example texts from “bild.de” and “faz.net” using a shellscript that downloaded the RSS-feeds of the newspaper-websites and fetched all its articles (~40 articles fetched) and stripped the relevant text out of the HTML-sourcecode.

All the sample texts were stored as several text-files in a directory for each language (“no/”, “sv/”) and newspaper (“bild/”, “faz/”) for later use (statistics!)

differential language features In order to find features that differ from language to language I wrote a script that analyzed the frequency of several things (most frequent words, most frequent characters and most frequent two-character-combinations). By looking at the data I concluded that the two-character-combinations could allow me to differentiate the languages from each other, because there was a visible difference in the order of the most frequent ones.

So I took the relative frequencies of the most frequent two-character-combinations and made a graph out of it (see Figure 1)



language	distinct character combinations	
sv	ar,at,fö,ll,oc,ra	(all green-only-bars)
no	fo,ha,le,ne,og,se	(all yellow-only-bars)

highest distinct bars: og (yellow) / oc (green)

Figure 1: relative frequencies of the most frequent two-character-combinations of Norwegian (yellow) and Swedish (green) texts

My idea how to differentiate the languages from each other was now to compare the relative

frequencies of the two-character-combinations which showed the highest deviation for particular languages in the graph.

So I wrote a initial R program that compared the frequency of “oc”, “og” (and “ei” for German, for the “extra challenge”) in a given file.

testing Then I wrote a script that called my inital R program several times with a portion of text of different length for each language and summed up the results.

I put these results into a graph (see Figure 2), then I tried to improve my R program using the results I got out of the graph and repeated the test run...

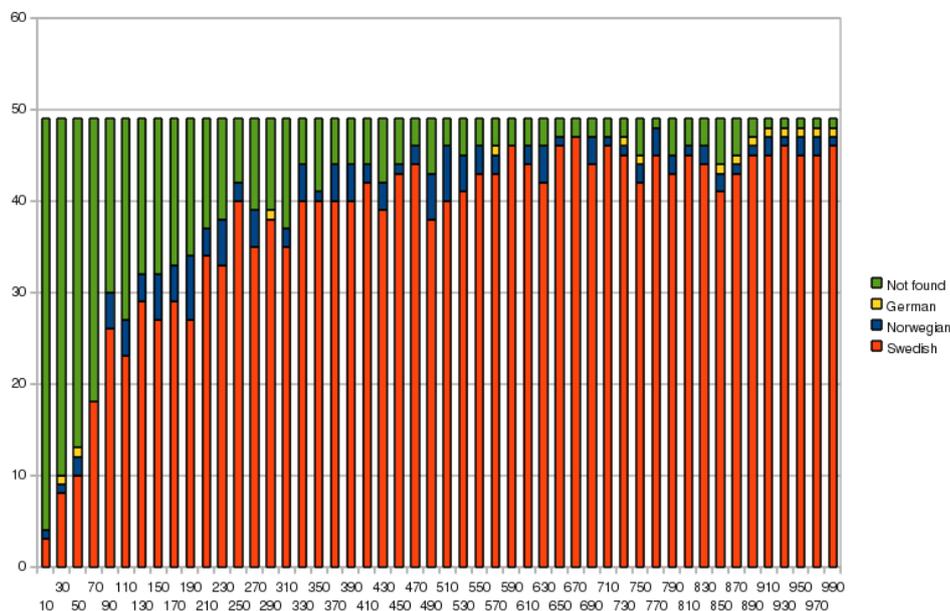


Figure 2: recognition results of the **initial** program with random Swedish (red) sample text portions with various lengths from 10 to 1000 characters

I repeated to improve and test the program until the graph showed a better recognition rate than the initial graph (see Figure 3): I adjusted the limits of the relative frequencies (i.e. how much percentage points are necessary for the relative frequency of a two-character-combination to identify a language) and I added a fallback routine which is called when the language couldn’t be identified in the first run. This routine is just like the main one, but it searches for other language-specific two-character-combinations (as you can see in Figure 1 there are more than just “og” and “oc”, in my case I’m using “ar” and “se”). This fallback routine increases the possibility for a short text to be identified, because especially in short texts there might be no occurrences of “og” and “oc” at all, and without this fallback it would simply not get identified, whereas with this fallback it has got a “second chance”.

error rates The estimated error rates are (roughly) apparent from the graph of Figure 3, but here are the values for the lengths 20, 200 and 2000 (as requested):

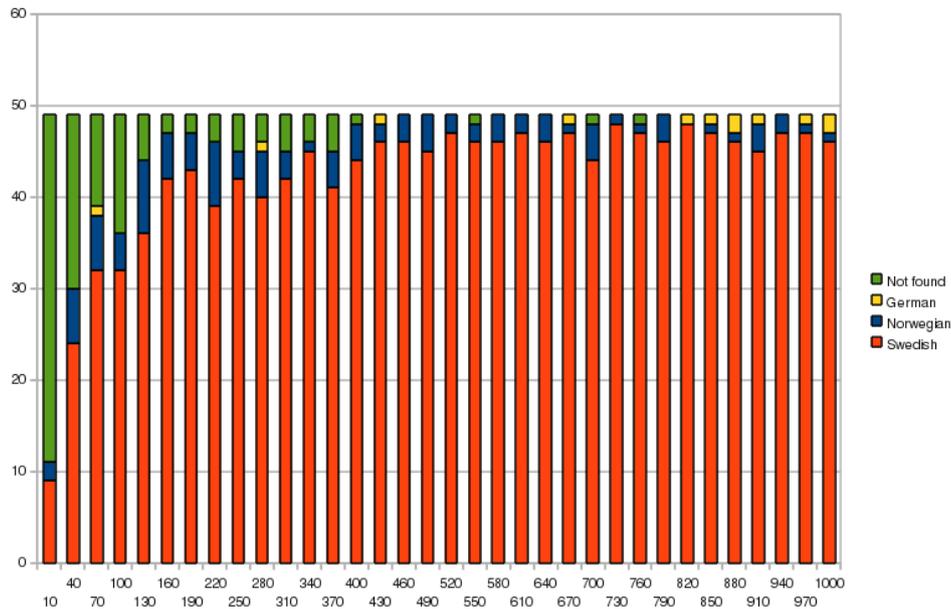


Figure 3: recognition results of the **improved** program with random Swedish (red) sample text portions with various lengths from 10 to 1000 characters

length	language	identified correctly	identified incorrectly	not identified	# tests
20	Swedish	134 (26%)	16 (3.2%)	349 (69.8%)	500
20	Norwegian	101 (20.2%)	52 (10.4%)	346 (69.2%)	500
200	Swedish	245 (81.6%)	40 (13.3%)	14 (4.6%)	300
200	Norwegian	236 (78.6%)	43 (14.3%)	20 (6.6%)	300
2000	Swedish	190 (95%)	9 (4.5%)	0 (0%)	200
2000	Norwegian	191 (95.5%)	8 (4%)	0 (0%)	200

limitations

- the algorithm uses statistical analysis which can never detect the language for sure, you can only assume that for a certain probability (depending on factors like length, type of text, ...) the computed result is right.
- the shorter the portion of text is the lower the probability of a correct language detection is (see graph Figure 3 or table error rates)
- if the text contains large amounts of non-language-elements (like names, formulas or other data) the detection is likely to be wrong because these non-language-elements can't be stripped out and will be used for language detection, which may falsify the result.

1.2 “Extra challenge“: distinguish texts taken from the FAZ from those you find in Bild

algorithm This algorithm is based on the assumption that a text with a high count of “?” and “!” is likely to be a “Bild”-text, because Bild quite often uses sensational expressions, in contrast to the FAZ, which uses a more complex and less emotional language.

error rates The error rates for distinguishing Bild from FAZ are worse than those for Swedish from Norwegian (see Figures 4 and 5). Much longer texts are needed to be able to identify either Bild or FAZ.

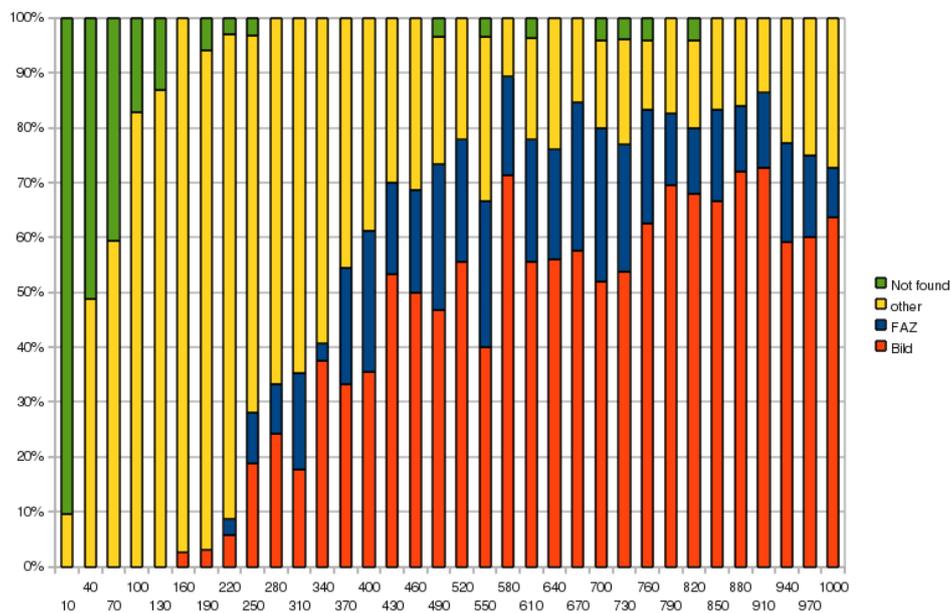


Figure 4: recognition results with random Bild (red) sample text portions with various lengths from 10 to 1000 characters

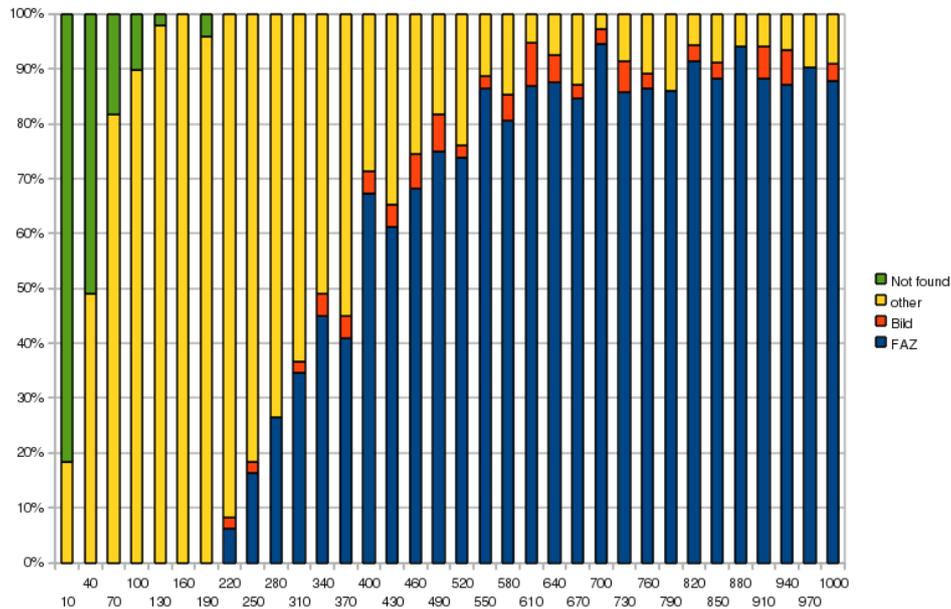


Figure 5: recognition results with random FAZ (blue) sample text portions with various lengths from 10 to 1000 characters

Listing 1: projekt1.R

```

1 #!/usr/bin/Rscript
2 # Genomik und Bioinformatik I
3 # Projekt 1 (Identify languages by text analysis)
4 # written by Andreas Loibl (enrollment no. 1524148)
5
6 source(file="readText.R")
7
8 # identifyGermanNewspaper:
9 #####
10 # distinguishes texts from the German newspapers FAZ and Bild by analyzing
11 # the relative frequency of punctuation marks in the text. This algorithm
12 # is based on the assumption a text with a high count of "?" and "!" is
13 # likely to be a "Bild"-text, because Bild quite often uses sensational
14 # expressions.
15 #
16 # So this function is able to decide if the given text is taken from Bild
17 # or from FAZ (or generic "Text" if the data is not distinguishable enough
18 )
19 identifyGermanNewspaper <- function(inputText) {
20   count <- numeric()
21   for(i in c("!", "\\?", "\\._", ",","")) {
22     tmp <- unlist(gregexpr(i, inputText))
23     if(tmp[1] > 0) count[i] <- length(tmp) else count[i] <- 0
24   }
25   if(sum(count) > 5) {
26     if((count["\\?"] + count["!"]) / sum(count) > 0.1) return("Bild")
27     if((count["\\?"] + count["!"]) / sum(count) < 0.05) return("FAZ")

```

```

28     print("Data_is_not_distinguishable_enough.")
29     return("Text")
30   } else {
31     print("No_or_not_enough_processable_data_found.")
32     return("Text")
33   }
34 }
35
36 # identifyLanguageOfFile:
37 #####
38 # distinguishes the language of a file by using text analysis on its
39 # content.
40 #
41 # the algorithm counts the occurrences of several character-combinations
42 # (like "oc", "og", "ei", "ar", "se") and estimates what language the text
43 # is written in by comparing the relative frequencys of these character-
44 # combinations with values.
45 identifyLanguageOfFile <- function(file) {
46   inputText <- readText(file)
47   count <- numeric()
48   for(i in c("og", "oc", "ei")) {
49     tmp <- unlist(gregexpr(i, inputText))
50     if(tmp[1] > 0) count[i] <- length(tmp) else count[i] <- 0
51   }
52   if(sum(count)>0) {
53     if((count["og"])/sum(count) > 0.40) return("Norwegian")
54     if((count["oc"])/sum(count) > 0.50) return("Swedish")
55     if((count["ei"])/sum(count) > 0.45) return(paste("German", "(",
56       identifyGermanNewspaper(inputText), ")"))
57   }
58   count <- numeric()
59   for(i in c("ar", "se")) {
60     tmp <- unlist(gregexpr(i, inputText))
61     if(tmp[1] > 0) count[i] <- length(tmp) else count[i] <- 0
62   }
63   if(sum(count)>0) {
64     if((count["ar"])/sum(count) > 0.65) return("Swedish")
65     if((count["se"])/sum(count) > 0.65) return("Norwegian")
66     print("Data_is_not_distinguishable_enough.")
67     return(FALSE)
68   }
69   print("No_processable_data_found.")
70   return(FALSE)
71 }
72 for(file in commandArgs(TRUE)) print(identifyLanguageOfFile(file))

```