



reddit

INSIGHTS FROM R/SKINCAREADDICTION

**Fetch reddit comments using package
Reddit ExtractoR**

Pre-process Text

```
#TEXT MINING
#using the tm package
library(tm)

#cleaning text as a part of preprocess, taking comments and putting into
#separate dataframe
commsposts <- sa_comms$comment

#dumping comments into corpus for preprocessing (using tm package)
sk_corpus <- Corpus(VectorSource(commsposts))

#change to lowercase
sk_corpus <- tm_map(sk_corpus, content_transformer(tolower))
#remove numbers
sk_corpus <- tm_map(sk_corpus, removeNumbers)
#remove punctuation
sk_corpus <- tm_map(sk_corpus, removePunctuation)
#remove stop words and "the" and "and"
sk_corpus <- tm_map(sk_corpus, removeWords, c("the", "and", stopwords("english")))
#stemming words
sk_corpus <- tm_map(sk_corpus, stemDocument, language = "english")
#removing white space
sk_corpus <- tm_map(sk_corpus, stripWhitespace)
```

TF-

IDF

=

$TF(t) = (\text{Number of times term appears in a comment}) / (\text{Total number of terms in the comment})$

*

$IDF(t) = \log(\text{Total number of comments} / \text{Number of comments with term } t \text{ in it})$

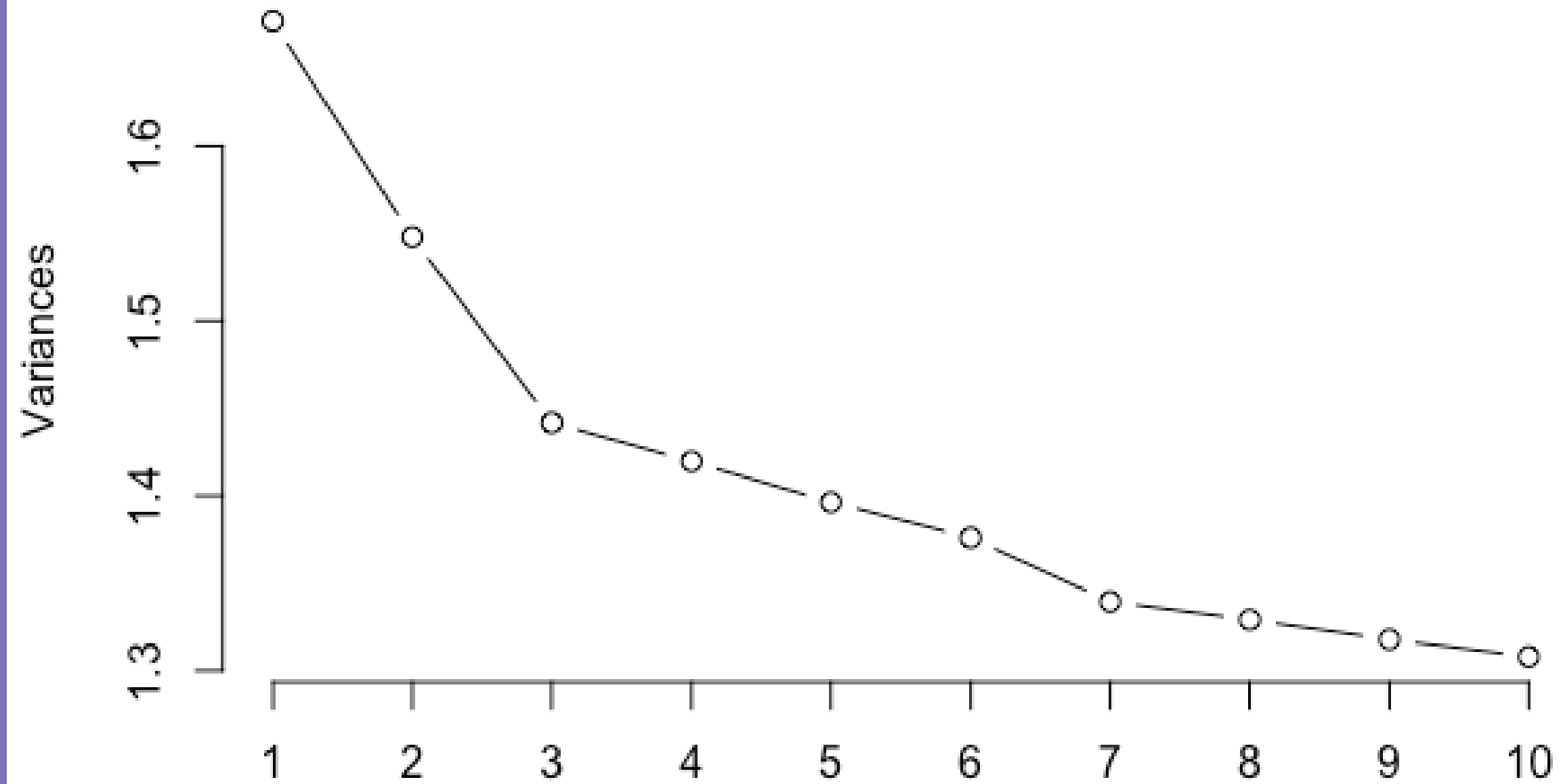
```
#ranking of terms
termsdec <- data.frame(sort(colSums(as.matrix(tfidfdf_skdf)), decreasing=TRUE))
termsdec #show decreasing list of tfidf words

library(wordcloud)
#now making wordcloud
wordcloud(rownames(termsdec), termsdec[,1], max.words=60, colors=brewer.pal(4, "Dark2"))
#word cloud shows that sunscreen, moisture, oil are high tf-idf terms
```

can time cream back
got thing might
also week see oil its
remov say find ive
still pretti spf never
peopl lot wash eye
get toner acn better
year appli even differ well
feel moistur sure

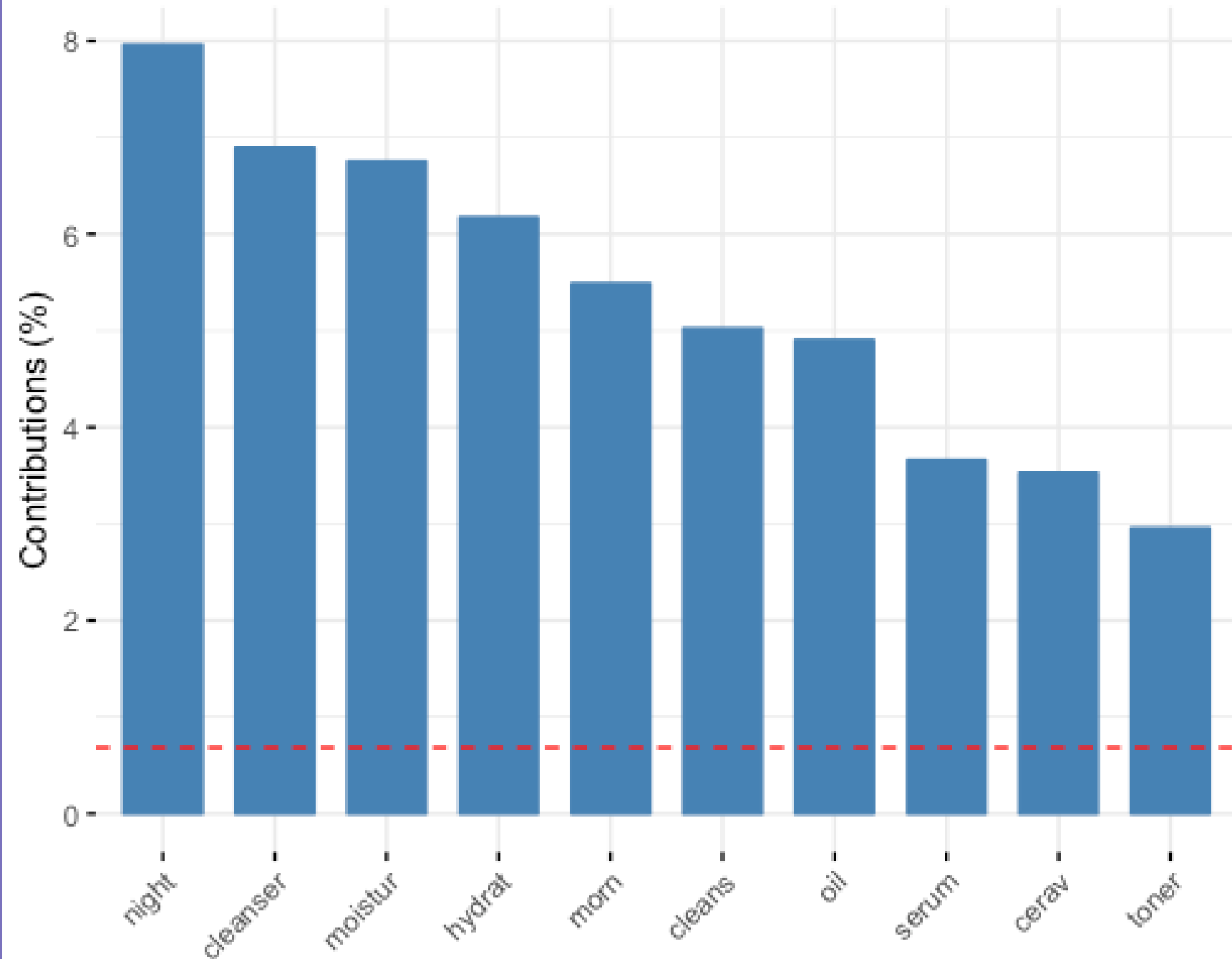
PCA

Scree Plot



PCA

Contribution of variables to Dim-1



Contribution of variables to Dim-2

