

Project 4 - Paper 5 Main Script

Xuehan Liu

04/05/2017

In this file, we illustrate our step-by-step procedure on the error driven online training algorithm on the nameset AKumar.txt (step0 - step3). Then we created a function called `error.driven.train()`, based on the step-by-step procedure described from step0-step3, and saved it in the “lib” folder. In step4, we implement the algorithm to all namesets provided and report the trained parameters.

Step 0: Load the packages, specify directories

```
setwd("/Users/xuehan/Desktop/Spr2017-proj4-team-9/")
# here replace it with your own path or manually set it in RStudio
# to where this rmd file is located

if (!require("pacman")) install.packages("pacman")

## Loading required package: pacman
pacman::p_load(text2vec, dplyr, qmcMatrix, kernlab, knitr)
```

Step 1: Load and process the data

```
AKumar <- data.frame(scan("../data/nameset/AKumar.txt",
                        what = list(Coauthor = "", Paper = "", Journal = ""),
                        sep=">", quiet=TRUE), stringsAsFactors=FALSE)

# extract canonical author id before "_"
AKumar$AuthorID <- sub("_.*", "", AKumar$Coauthor)
# extract paper number under same author between "_" and first whitespace
AKumar$PaperNO <- sub(".*_(\\w*)\\s.*", "\\1", AKumar$Coauthor)
# delete "<" in AKumar$Coauthor, you may need to further process the coauthor
# term depending on the method you are using
AKumar$Coauthor <- gsub("<", "", sub("^.*?\\s", "", AKumar$Coauthor))
# delete "<" in AKumar$Paper
AKumar$Paper <- gsub("<", "", AKumar$Paper)
# add PaperID for further use, you may want to combine all the nameset files and
# then assign the unique ID for all the citations
AKumar$PaperID <- rownames(AKumar)
```

Step 2: Feature Design

As mentioned in the paper, we can use TF-IDF to collect all unique terms in each citation.

```
it_train <- itoken(AKumar$Paper,
                  preprocessor = tolower,
                  tokenizer = word_tokenizer,
```

```

ids = AKumar$PaperID,
# turn off progressbar because it won't look nice in rmd
progressbar = FALSE)
vocab <- create_vocabulary(it_train, stopwords = c("a", "an", "the", "in", "on",
"at", "of", "above", "under"))

#vocab

vectorizer <- vocab_vectorizer(vocab)
dtm_train <- create_dtm(it_train, vectorizer)
dim(dtm_train)

## [1] 244 666

tfidf <- TfIdf$new()
dtm_train_tfidf <- fit_transform(dtm_train, tfidf)

```

Step 3: Implementing hierarchical clustering and training parameters

In this section, our goal is to train the lambda on hierarchical clustering on our text file AKumar by the error driven online training method introduced in the paper5. We use the ranking perceptron to update the parameters.

```

####Initialize Parameter lambda
lambda<-rep(1,nrow(dtm_train_tfidf))

#Add the Author's ID as the label column to the feature matrix for future use
dtm_train_tfidf<-cbind(dtm_train_tfidf,as.numeric(AKumar$AuthorID))

#Given the training set, we are able to generate the true clusters.
#Based on the paper, we define true score S_star as the distance of the sum of clusterwise distance. We

#Compute the true score S_star for the giving training data
element<-list()
S_star<-vector(length=length(unique(AKumar$AuthorID)))
for (i in 1:length(unique(AKumar$AuthorID))){
  element[[i]]<-dtm_train_tfidf[dtm_train_tfidf[,ncol(dtm_train_tfidf)]==i,]
  S_star[i]<-sum(dist(element[[i]]))/2
}
S_star<-mean(S_star)
T_star<-dtm_train_tfidf[,ncol(dtm_train_tfidf)]

K=14
k=1
lambda1 <- matrix(NA, nrow = nrow(AKumar), ncol = K)
S1 <- numeric(K)
acc <- numeric(K)
while (k<(K+1)){

#Implement Hierarchical Clustering
h<-hclust(dist(dtm_train_tfidf*lambda))
#Check the result for the number of cluster equals to the number of unique authors in the dataset.

```

```

h_result<-cutree(h,k=length(unique(AKumar$AuthorID)))

#Compute the our own score function S
S<-vector(length=length(unique(AKumar$AuthorID)))
element_s<-list()
for (i in 1:length(unique(AKumar$AuthorID))){
  element_s[[i]]<-dtm_train_tfidf[which(h_result==i),]
  S[i]<-sum(dist(element_s[[i]]))/2
}
S<-mean(S)

#Identify true author for each cluster generated by hclust() function, and assign it to each element of

label<-dtm_train_tfidf[,ncol(dtm_train_tfidf)]
author.clust <- vector(length=length(unique(AKumar$AuthorID)))
for (i in 1:length(unique(AKumar$AuthorID))){
  author.clust[i]<-as.numeric(names(which.max(table(label[which(h_result==i)]))))
}

for (i in 1:unique(AKumar$AuthorID)){
  h_result[h_result==i]<-author.clust[i]
}
T_hat<-h_result

#Update lambda

for (i in 1:length(T_star)){
  if (T_hat[i]!=T_star[i]){
    lambda[i]<-lambda[i]-((S-S_star)/S) ###
  }
  else {
    lambda[i]<-lambda[i]
  }
}
lambda1[,k] <- lambda
S1[k] <- S
acc[k] <- mean(label == h_result)
k=k+1
}

```

```

## Warning in 1:unique(AKumar$AuthorID): numerical expression has 14 elements:
## only the first used

```

```

## Warning in 1:unique(AKumar$AuthorID): numerical expression has 14 elements:
## only the first used

```

```

## Warning in 1:unique(AKumar$AuthorID): numerical expression has 14 elements:
## only the first used

```

```

## Warning in 1:unique(AKumar$AuthorID): numerical expression has 14 elements:
## only the first used

```



```
## [1] 0.9404979
```