

# Paper5 (C/E/Mr)

Yu Shan & Ke HAN

## Step 0: Load packages, specify directories

```
#setwd(".....")  
# which would help with the df csv file saving
```

## Step 1: Load and process the data

```
# final version of data processing  
# similar with NB_CO  
# do not paste here, using same data_list  
  
data_gather=function(){  
  aa=1:14  
  temp=c("df1.csv","df2.csv","df3.csv","df4.csv","df5.csv","df6.csv","df7.csv",  
    "df8.csv","df9.csv","df10.csv","df11.csv","df12.csv","df13.csv","df14.csv")  
  for(hh in aa){  
    hh=1:14  
    num.paper<-length(data_list[[hh]])  
  
    author.id<-NULL  
    coauthor<-NULL  
    for(i in 1:num.paper){  
      author.id[i]<-data_list[[hh]][[i]][[1]]  
      coauthor<-unique(c(coauthor,data_list[[hh]][[i]][[3]]))  
    }  
  
    df<-data.frame(author.id)  
  
    num.coauthor<-length(coauthor)  
  
    df<-cbind(author.id,id_co)  
  
    get.coauthor<-function(i){  
      coauthors<-rep(0,num.coauthor)  
      for(j in 1:length(data_list[[hh]][[i]][[3]])){  
        coauthors[which(data_list[[hh]][[i]][[3]][j]==coauthor)]<- 1  
      }  
      return(coauthors)  
    }  
    id_co<-NULL  
    for(i in 1:num.paper){  
      id_co<-rbind(id_co,get.coauthor(i))  
    }  
  }  
}
```

```

}
colnames(id_co)<-coauther
df<-cbind(auther.id,id_co)
write.csv(df,temp[hh])
cat(hh)
cat('....')
cat(' ')
}
}

```

## Step 2: C/ E/ Mr

*# Clusterwise Scoring Function & Error-driven Online Training & Ranking  
MIRA*

```

#final function
update.weights=function(data,weights,target)
{
  pred=which.max(weights%%t(data.matrix(data)))
  tau=(t(data.matrix(weights[pred,]-weights[target,]))%%t(data.matrix(data)))+1)/(2*sum(data^2))
  if(pred!=target)
  {
    weights[pred,]=as.vector(weights[pred,])-as.vector(unlist(min(tau[1,1],0.008)*data))
    weights[target,]=as.vector(weights[target,])+as.vector(unlist(min(tau[1,1],0.008)*data))
  }

  return(weights)
}

predict.mira=function(data,weights)
{
  return(apply(weights%%t(data.matrix(data)),2,which.max))
}

mira=function(x,y,levels=length(unique(y)))
{
  #y=as.numeric(as.factor(y))
  weights=matrix(rep(1,levels*length(x[1,])),nrow=levels)
  weights=update.weights(x[1,],weights,y[1])
  pred=predict.mira(x,weights)
  errorid=which(pred!=y)
  diff=1

```

```

while(diff>0 && (length(errorid)!=0))
{
weights=update.weights(x[errorid[1],],weights,y[errorid[1]])
pred=predict.mira(x,weights)
errorid2=which(pred!=y)
diff=length(errorid)-length(errorid2)
errorid=errorid2
}
return(weights)
}

```

### Step 3: Evaluation

Calculate the accuracy of the 14 authors.

```

run=function(){
  res=c()
  #set wd to get dfs

  temp=c("df1.csv","df2.csv","df3.csv","df4.csv","df5.csv","df6.csv","d
f7.csv","df8.csv","df9.csv","df10.csv","df11.csv","df12.csv","df13.csv",
"df14.csv")
  for (i in 1:length(temp)) {
    df = read.csv(temp[i], header = TRUE)
    x=df[,-c(1,2)]
    y=as.numeric(as.factor(df$auther.id))
    trainid=sample(1:length(y),length(y)/2)
    trainx=x[trainid,]
    trainy=y[trainid]
    testx=x[-trainid,]
    testy=y[-trainid]
    weights=mira(trainx,trainy,levels=length(unique(y)))
    o=predict.mira(testx,weights)
    res[i]=sum(testy==o)/length(testy)
    cat(i)
    cat(' ')
  }
  return (res)
}
run()

## 1 2 3 4 5 6 7 8 9 10 11 12 13 14

## [1] 0.45674740 0.84426230 0.16708229 0.91304348 0.06338028 0.875000
00
## [7] 0.82558140 0.30172414 0.87142857 0.79220779 0.85384615 0.635922
33
## [13] 0.09836066 0.42812006

```