

1 Topic modelling with NMF

```
Ptilde <- as.matrix( read.csv("data/news.csv") )  
P <- Ptilde/sum(Ptilde)  
showcol(pmin(P,0.0001))
```

a)

```
r <- 4  
lr.gkl <- lee01.gkl(P, r, reps=5)  
with(lr.gkl,  
  for (k in 1:nrow(R)) {  
    print(rev(sort(R[k,]))[1:10])  
    cat(strrep('-',130), "\n")  
  })
```

```
ggplotm(lr.gkl$R[,1:10], format="", show.axis=FALSE, mid="black")
```

b)

```
Phat <- lr.gkl$L %*% lr.gkl$R  
showcol(pmin(Phat,0.0001))
```

c)

ii)

```
Phat <- P.svd$u[,1:4] %*% diag(P.svd$d[1:4]) %*% t(P.svd$v[,1:4])  
showcol(pmin(Phat,0.0001))
```

```
r <- 2  
lr.gkl <- lee01.gkl(P, r, reps=5)  
with(lr.gkl,  
  for (k in 1:nrow(R)) {  
    print(rev(sort(R[k,]))[1:10])  
    cat(strrep('-',130), "\n")  
  })
```

```
Phat <- lr.gkl$L %*% lr.gkl$R  
showcol(pmin(Phat,0.0001))
```

ii)

```
r <- 8
lr.gkl <- lee01.gkl(P, r, reps=5)
with(lr.gkl,
  for (k in 1:nrow(R)) {
    print(rev(sort(R[k,]))[1:10])
    cat(strrep('-',130), "\n")
  })
```

```
Phat <- lr.gkl$L %*% lr.gkl$R
showcol(pmin(Phat,0.0001))
```

e)

```
r <- 4
lr.gnmf <- lee01.gnmf(P, r, reps=5)
with(lr.gnmf,
  for (k in 1:nrow(R)) {
    print(rev(sort(R[k,]))[1:10])
    cat(strrep('-',130), "\n")
  })
```

```
Phat <- lr.gnmf$L %*% lr.gnmf$R
showcol(pmin(Phat,0.0001))
```

2 PLSA

a)

```
lsr.gkl <- nmf.lsr(lr.gkl)
summary(lsr.gkl)
```

	Length	Class	Mode
L	1600	-none-	numeric
S	16	ddiMatrix	S4
R	3200	-none-	numeric

```
> apply(lsr.gk$L,2,sum)
[1] 1 1 1 1
> sum(lsr.gk$S)
[1] 0.999818
> apply(lsr.gk$R,1,sum)
[1] 1 1 1 1
```

b)

```
slr.gkl <- nmf.slr(lr.gkl)
summary(slr.gkl)
```

	Length	Class	Mode
S	160000	ddiMatrix	S4
L	1600	-none-	numeric
R	3200	-none-	numeric

[illegible]

```
[388] 1 1 1 1 1 1 1 1 1 1 1 1 1 1  
> apply(slr.gkl$R,1,sum)  
[1] 1 1 1 1
```

3 Clustering

a)

```
cluster <- kmeans(P, 4, nstart=100)$cluster
```

```
> cluster
```

[illegible]

Relable:

```
i<-1
while(i<101){
  if(cluster[i]==1){cluster[i]<-4
}
else if(cluster[i]==2){cluster[i]<-2
}
else if(cluster[i]==3){cluster[i]<-3
}
else{cluster[i]<-1
}
i<-i+1
}
trail_a<- cm(cluster)

trail_a$overall["Accuracy"]
```

b)

```
P.svd <- svd(P)
```

```
test2<- P.svd$u[,1:4]%*%diag(P.svd$d[1:4])
```

```
cluster <- kmeans(test2, 4, nstart=100)$cluster
```

```
> cluster
```

[illegible]

Relabel:

```
i<-1
while(i<101){
  if(cluster[i]==1){cluster[i]<-4
}
else if(cluster[i]==2){cluster[i]<-2
}
else if(cluster[i]==3){cluster[i]<-3
}
else{cluster[i]<-1
}
i<-i+1
}
```

```
trail_b<- cm(cluster)
```

```
trail_b$overall["Accuracy"]
```

c)

```
r <- 4
```

```
l.r.gk1 <- lee01.gk1(P, r, reps=5)
```

```
cluster <- kmeans(lr.gkl$L, 4, nstart=100)$cluster
```

```
> cluster
```

[illegible]

[321] 2 2 2 1 2 1 2 2 2 2 2 1 2 2 2 1 1 2 2 2 2 2 1 1 1 2 2 1 2 1 1 2 2 2 1 2 2 1 1 2 1 2 2 2 1 2 2 2 2 2 2 2 1 2 2 2 2 1 2 2 2 2 2

Relabel:

```
i<-1
while(i<101){
  if(c1uster[i]==1){c1uster[i]<-4
}
else if(c1uster[i]==2){c1uster[i]<-2
}
else if(c1uster[i]==3){c1uster[i]<-3
}
else{c1uster[i]<-1
}
i<-i+1
}
```

```
trail_c<- cm(cluster)
```

```
trail_c$overall["Accuracy"]
```

d)

```
l1sr.gkl <- nmf.l1sr(l1r.gkl)
```

```
cluster <- kmeans(lsr.gkl$L, 4, nstart=100)$cluster
```

```
> cluster
```

[illegible]

Relabel:

```
i<-1
while(i<101){
  if(cluster[i]==1){cluster[i]<-3
}
else if(cluster[i]==2){cluster[i]<-2
}
else if(cluster[i]==3){cluster[i]<-4
}
else{cluster[i]<-1
}
i<-i+1
}
```

```
trail_d<- cm(cluster)
```

```
trail_d$overall["Accuracy"]
```

e)

```
s1r.gkl <- nmf.s1r(lr.gkl)
```

```
cluster <- kmeans(slr.gkl$L, 4, nstart=100)$cluster
```

```
> cluster
```

[illegible]

Relabel:

```
i<-1
while(i<101){
  if(cluster[i]==1){cluster[i]<-2
} else if(cluster[i]==2){cluster[i]<-1
} else if(cluster[i]==3){cluster[i]<-4
```



```
    }else{cluster[i]<-3  
    }  
    i<-i+1  
  }  
}
```

```
trail_e<- cm(cluster)
```

```
trail_e$overall["Accuracy"]
```