

CCA on Halper & Yang

Nick Wawee

August 27, 2018

Loading

```
#Loading
loadmaty<-read.table("Yang_ReadCount_ERCC_Converted.txt", fill = TRUE, stringsAsFactors = FALSE)
loadmaty2<-read.table("pseudo10mat.txt", stringsAsFactors = FALSE)
zonatedmat<-read.table("OrderedMatrix_Filtered_1.txt", stringsAsFactors = FALSE)
zonatedmat<-as.data.frame(zonatedmat)

#Feature Data
zgenes<-rownames(zonatedmat)
zonatedmat<-cbind(zgenes,zonatedmat)
zonatedmat$genes<-as.character(zonatedmat$genes)
zonatedmat<-zonatedmat %>%
  separate_rows(genes, sep=";")
rownames(zonatedmat2)<-zonatedmat2[,1]

# Testing separate rows
testexp<-zonatedmat2[,c(which(rownames(zonatedmat2)=="Apr3"),which(rownames(zonatedmat2)=="0610007C21R1k")),]
testexp2<-zonatedmat2[which(rownames(zonatedmat2)=="0610007C21R1k;Apr3"),,]
#View(rbind(testexp,testexp2))

#-----

loadmaty[,2]<-as.character(loadmaty[,2])
matching<-match(rownames(zonatedmat2),loadmaty[,2])
matching<-matching[which(is.na(matching))]
fdmat<-as.matrix(loadmaty[matching,1:2])
rownames(fdmat)<-fdmat[,1]
fdmat<-as.matrix(fdmat[,1])
colnames(fdmat)<-"gene_abund_name"
fd<-as.data.frame(fdmat)
fd<-new("AnnotatedDataFrame", data=fd)

#Pheno Data
#Yang
cellnamesy<-as.matrix(colnames(loadmaty2)[which(loadmaty2[,2,]== "Hepatoblast"|loadmaty2[,2,]== "Hepatocyte")])
colnames(cellnamesy)<-cellnamesy[,1,]
cellnamesy<-cbind(cellnamesy, rep("Zonation",length(cellnamesy[,1])))
colnames(cellnamesy)<-c("Cell_Name", "Cell_Type")
rownames(cellnamesy)<-cellnamesy[,1,]
pt<-loadmaty2[,match(rownames(cellnamesy),colnames(loadmaty2))]
pt<-as.numeric(pt)
pt<-scalefun(pt, "0and1")
statey<-loadmaty2[,match(rownames(cellnamesy),colnames(loadmaty2))]
statey<-t(statey)
colnames(statey)<- "State0"
statey2<-statey
colnames(statey2)<- "State2"
cellnamesy<-cbind(cellnamesy,pt,statey,statey2)

#Sonation
cellnamesz<-as.matrix(colnames(zonatedmat2)[2:length(zonatedmat2)])
rownames(cellnamesz)<-cellnamesz[,1,]
cellnamesz<-cbind(cellnamesz, rep("Zonation",length(cellnamesz[,1])))
colnames(cellnamesz)<-c("Cell_Name", "Cell_Type")
#cellnamesz<-cellnamesz[-length(cellnamesz[,1]),]
cellnamesz<-cbind(cellnamesz,t(zonatedmat2[c(1,4,5),-1]))
cellnamesz[,3]<-scalefun(as.numeric(colnamesz[,3])), "0and1")
colnames(cellnamesz)<-colnames(cellnamesy)
#Combination
pdmat<-as.matrix(rbind(cellnamesy,cellnamesz))
pddf<-as.data.frame(pdmat)
pddf$pt<-as.numeric(as.character(pddf$pt))
pd<-new("AnnotatedDataFrame", data=pddf)

#Numeric Matrices
#Yang
rownames(loadmaty)<-loadmaty[,1,]
maty<-loadmaty[2:length(loadmaty[,1]),match(rownames(cellnamesy),loadmaty[,1,])]
colnames(maty)<-rownames(cellnamesy)
maty<-as.matrix(maty)
matching<-match(rownames(fdmat),rownames(maty))
maty<-maty[matching,]
nummaty<-mat2numericmat(maty)
#idandname<-function(idanemat, geneid, genename)
rownames(nummaty)<-idandname(fdmat, rownames(nummaty), "NA")

## Warning in if (geneid == "NA") {: the condition has length > 1 and only the
## first element will be used

#Sonation
tidmatz<-zonatedmat2[6:length(zonatedmat2[,1]),]
matching<-match(rownames(tidmatz),fdmat[,1,])
matching<-matching[which(is.na(matching))]
tidmatz<-tidmatz[matching,]
rownames(tidmatz)<-rownames(fdmat)[matching]
matz<-as.matrix(tidmatz[,1])
nummatz<-mat2numericmat(matz)
rownames(nummatz)<-idandname(fdmat, rownames(nummatz), "NA")

## Warning in if (geneid == "NA") {: the condition has length > 1 and only the
## first element will be used
```

Seurat Implementation

```
#Yang cell dataset
cdsy <- CreateSeuratObject(raw.data = nummaty)
cdsy <- NormalizeData(object = cdsy)
cdsy <- FindVariableGenes(object = cdsy, do.plot = FALSE)
#Zonation cell dataset
cdsz <- CreateSeuratObject(raw.data = nummatz)
cdsz <- NormalizeData(object = cdsz)
cdsz <- ScaleData(object = cdsz)
cdsz <- FindVariableGenes(object = cdsz, do.plot = FALSE)

# we will take the union of the top 2k variable genes in each dataset for
# alignment note that we use 2k genes in the manuscript example, you can
# try this here with negligible changes to the overall results
hvg.cdsy <- rownames(x = head(x = cdsy$hvg.info, n = 4437))
hvg.cdsz <- rownames(x = head(x = cdsz$hvg.info, n = 4437))
hvg.union <- union(x = hvg.cdsy, y = hvg.cdsz)

cds$meta.data[, "protocol"] <- "Zonation"
cds$meta.data[, "protocol"] <- "Yang"
```

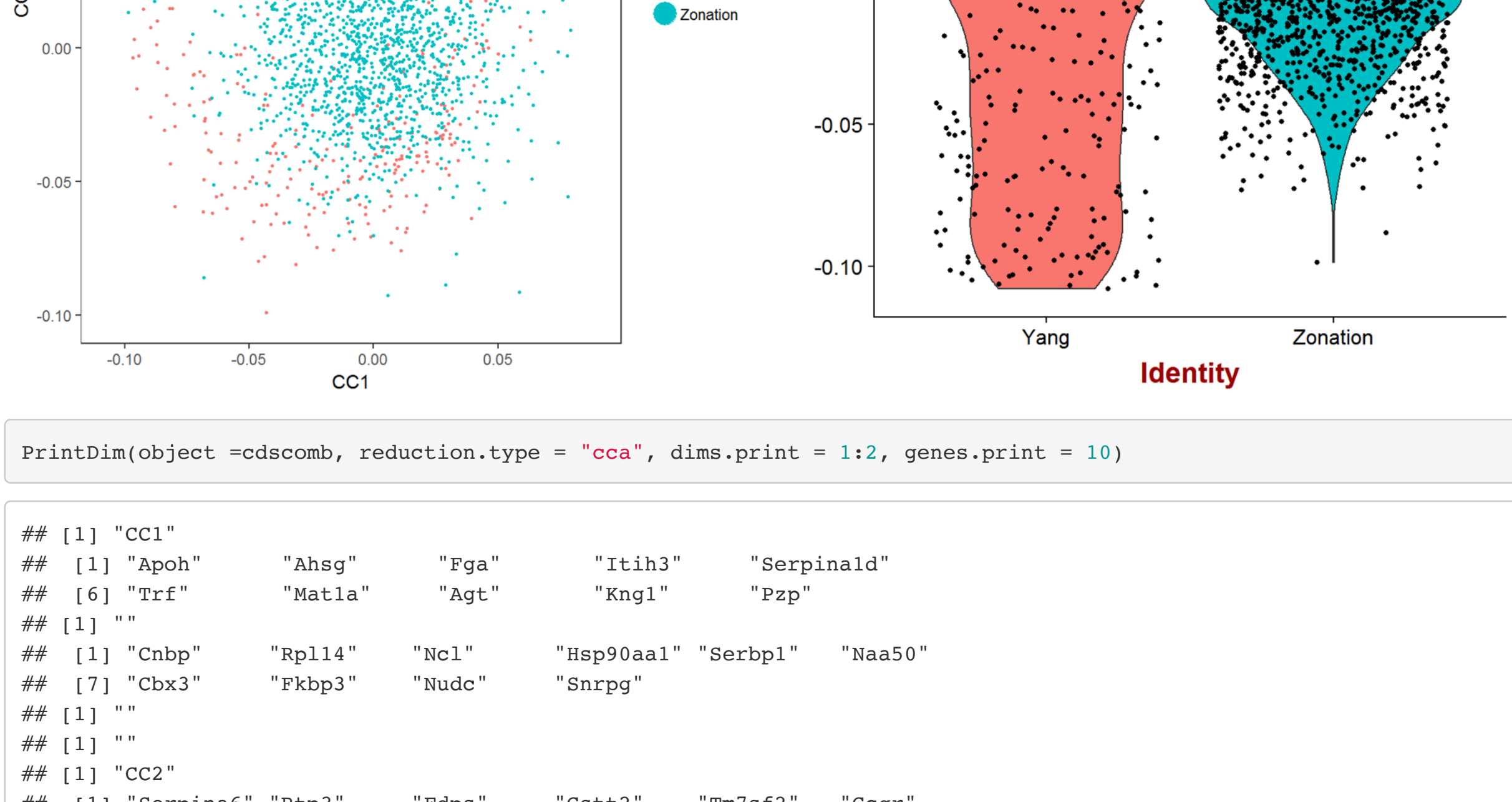
Running CCA

```
cdscomb<- RunCCA(object = cdsy, object2 = cdsz, genes.use = hvg.union)

## Warning: package 'bindrcpp' was built under R version 3.4.4

## Scaling data matrix

p1 <- DimPlot(object = cdscomb, reduction.use = "cca", group.by = "protocol", pt.size = 0.5,
do.return = TRUE)
p2 <- VlnPlot(object = cdscomb, features.plot = "CC1", group.by = "protocol", do.return = TRUE)
plot_grid(p1, p2)
```

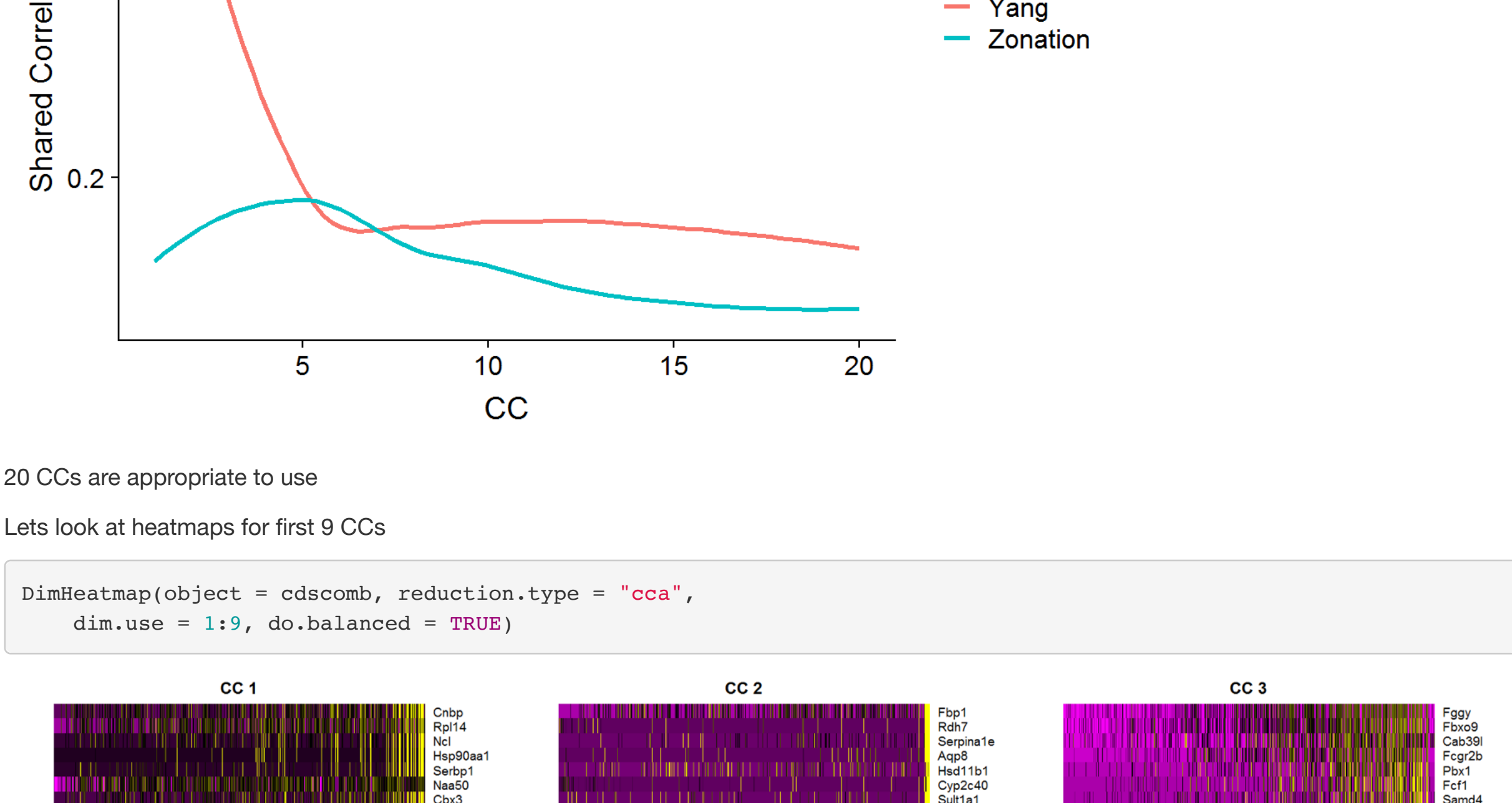


```
PrintDim(object =cdscomb, reduction.type = "cca", dims.print = 1:2, genes.print = 10)

## [1] "CC1"
## [1] "Aph" "Ahsg" "Fga" "Itih3" "Serpina1d"
## [6] "Trf" "Matla" "Kng1" "Kng1" "Pzp"
## [11] ""
## [1] "Ccbp" "Rpl14" "Ncl" "Hsp90aa1" "Serpi1" "Haa50"
## [7] "Cbx3" "Fkbp3" "Nudo" "Snpg"
## [13] ""
## [1] "CC2"
## [1] "Serpina4" "Rtp3" "Pdps" "Gstt2" "Tm7sf2" "Gcgr"
## [7] "Apoc2" "Rgtz34" "Cpt2" "Aldh3a2"
## [13] ""
## [1] "Fbp1" "Rdh7" "Serpina1e" "Aqp8" "Hsd11b1"
## [6] "Cyp2c40" "Sult1a1" "Serpina1c" "Qsox1" "Hmg1"
## [13] ""
## [1] ""
```

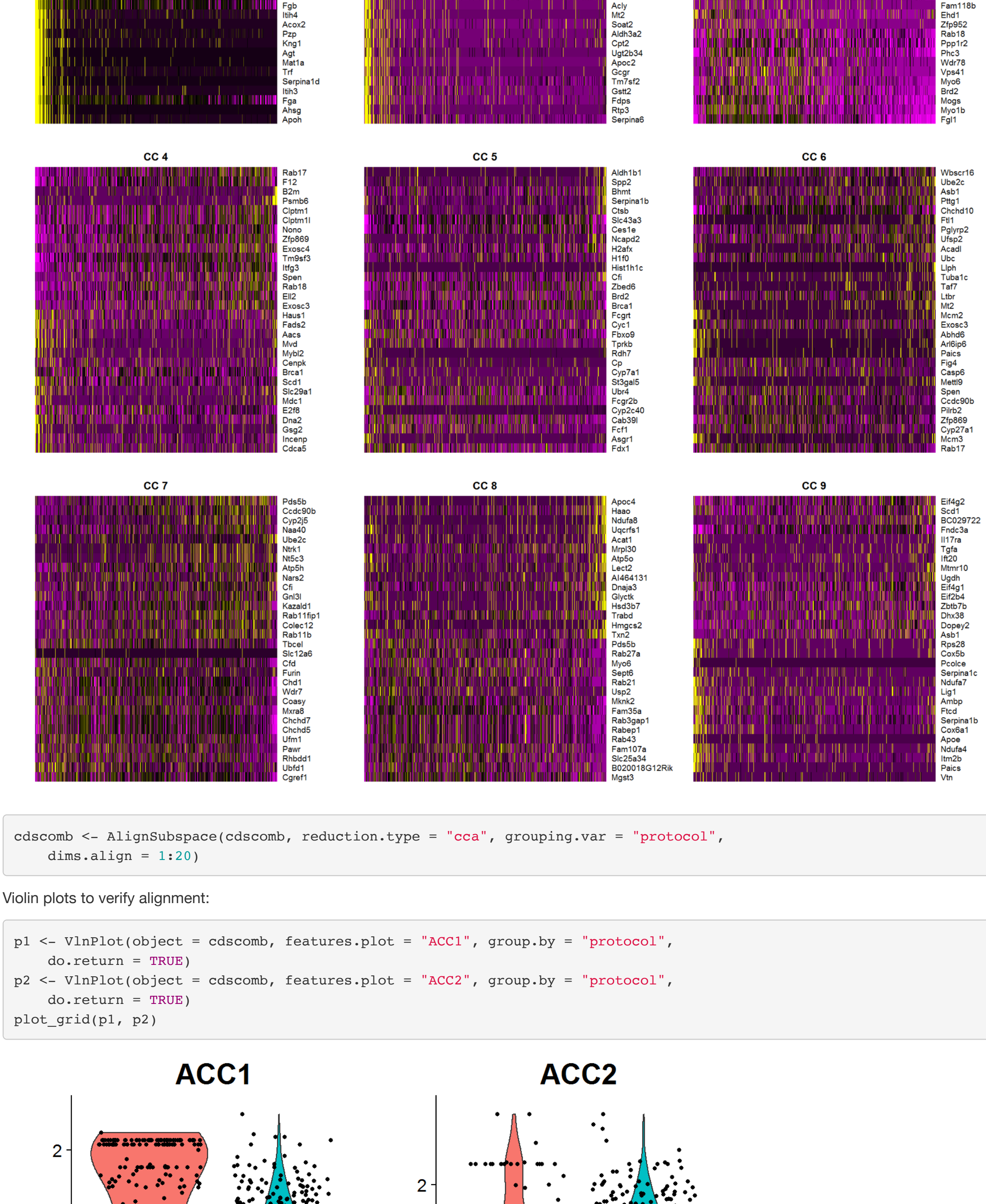
Choosing which CCs to use

```
p3 <- MetaGeneBicorPlot(cdscomb, grouping.var = "protocol", dims.eval = 1:20,
display.progress = FALSE)
```



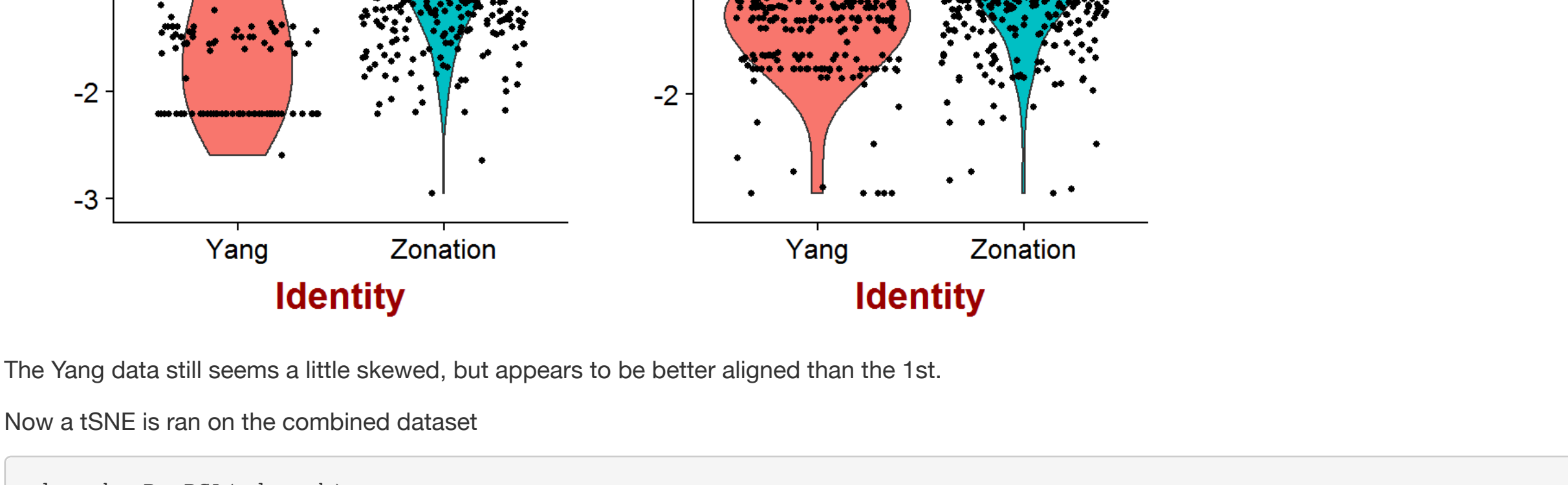
20 CCs are appropriate to use

Lets look at heatmaps for first 9 CCs



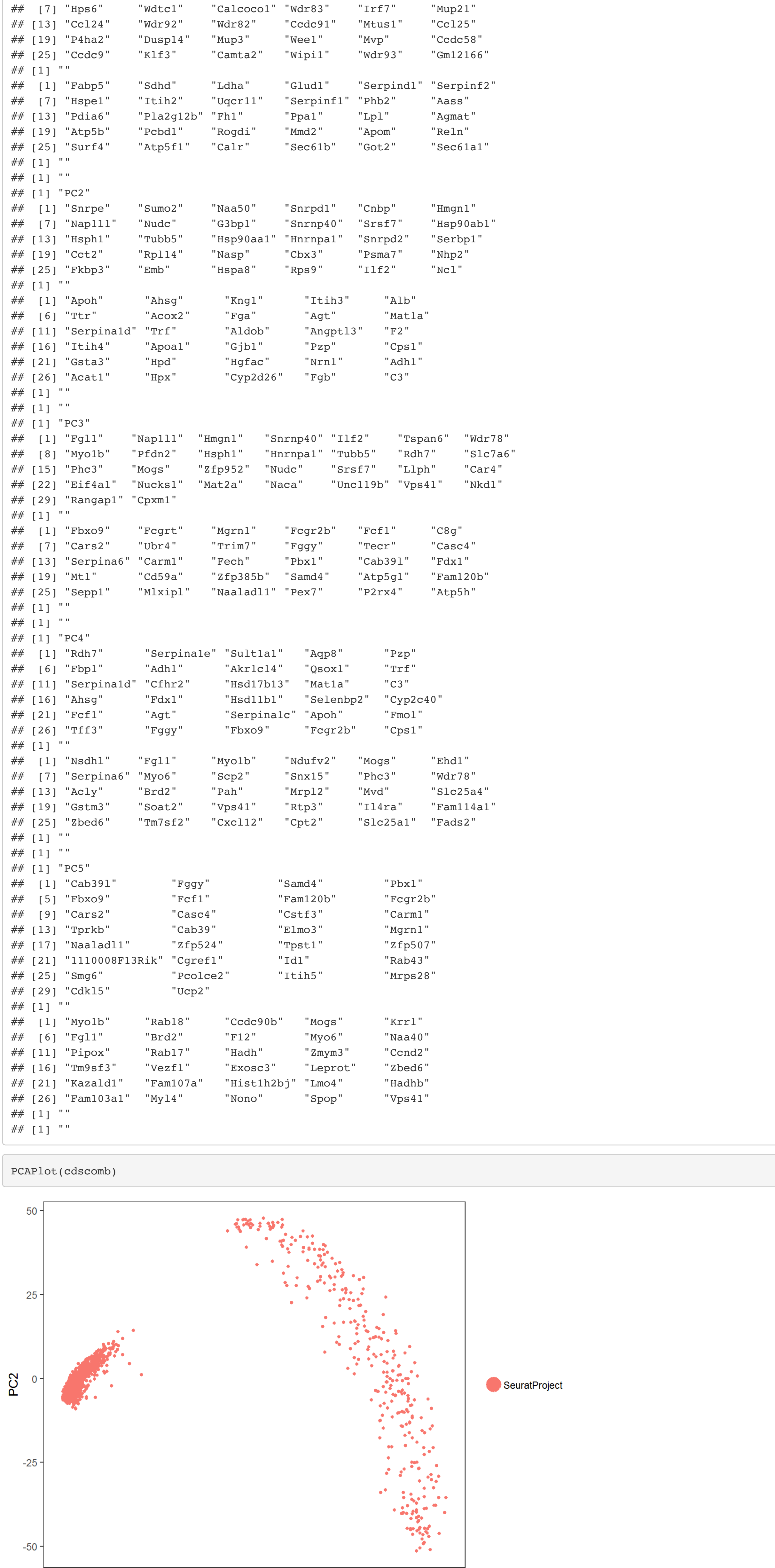
```
cdscomb <- AlignSubspace(cdscomb, reduction.type = "cca", grouping.var = "protocol",
dims.align = 1:20)
```

Violin plots to verify alignment:



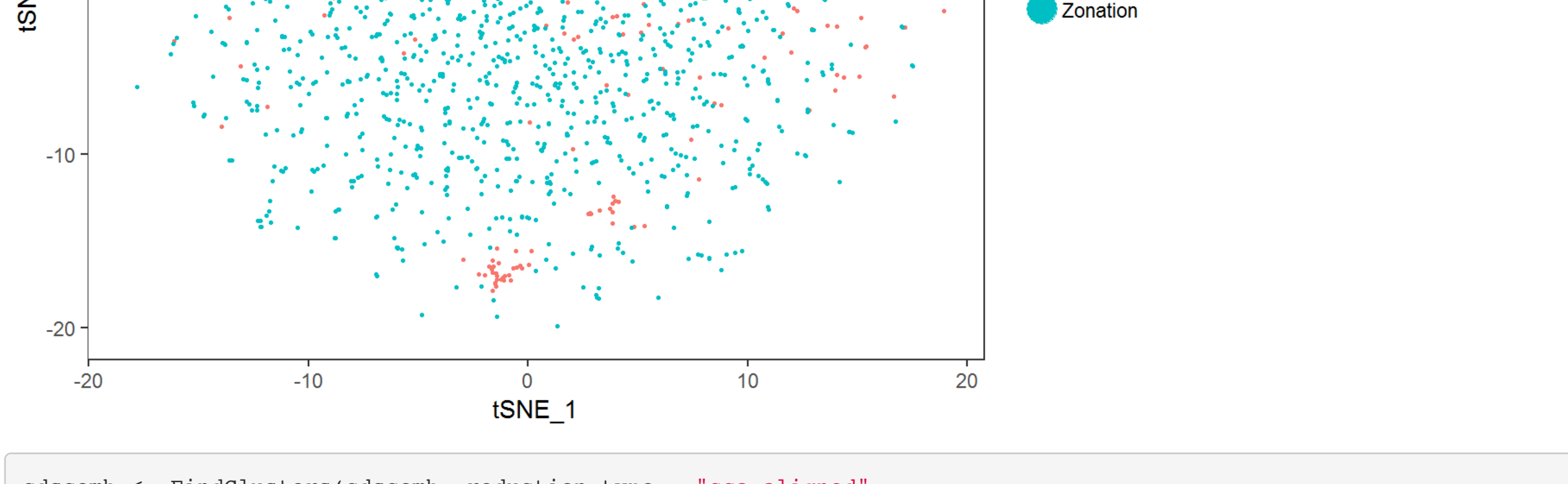
The Yang data still seems a little skewed, but appears to be better aligned than the 1st.

Now a tSNE is ran on the combined dataset

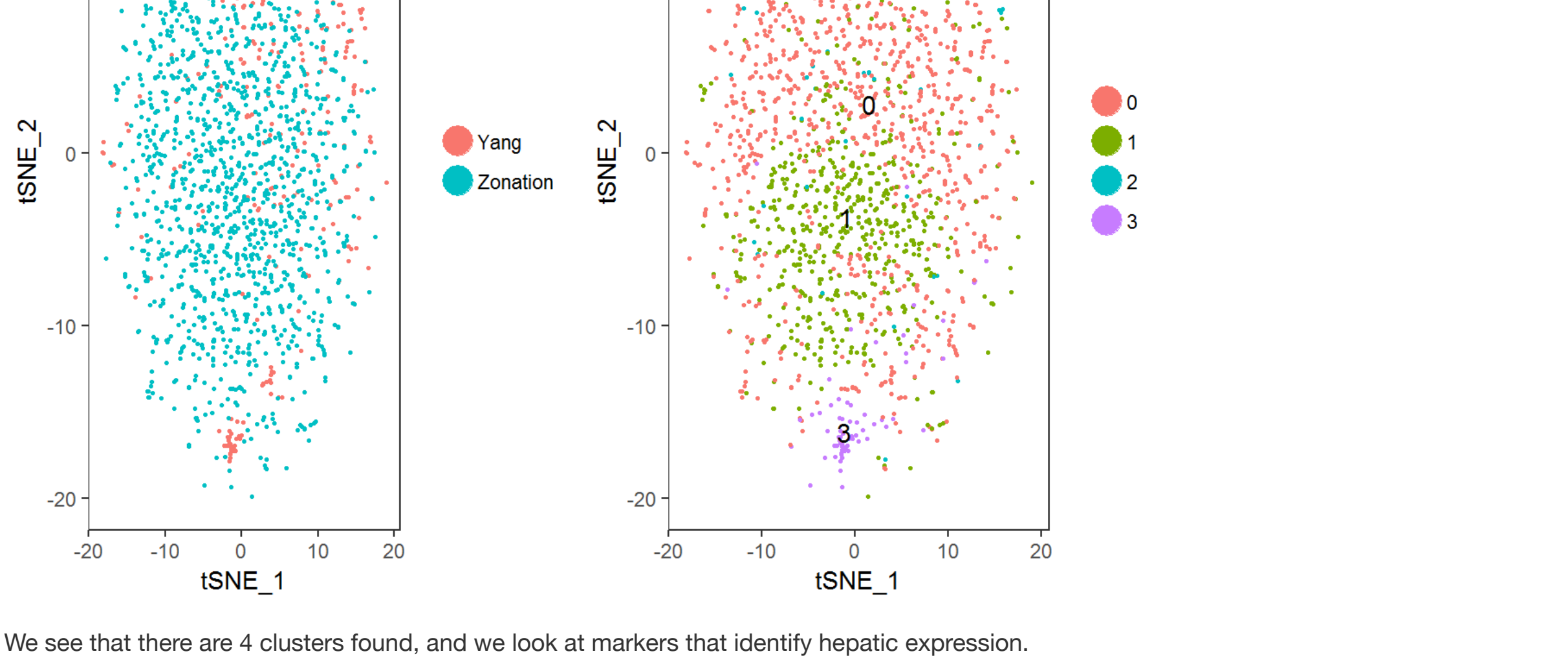


```
cdscomb<- RunTSNE(cdscomb, reduction.use = "cca,aligned", dims.use = 1:20,
do.fast = F)
```

```
p1 <- TSNEPlot(cdscomb, do.return = T, pt.size = 0.5, group.by = "protocol")
print(p1)
```

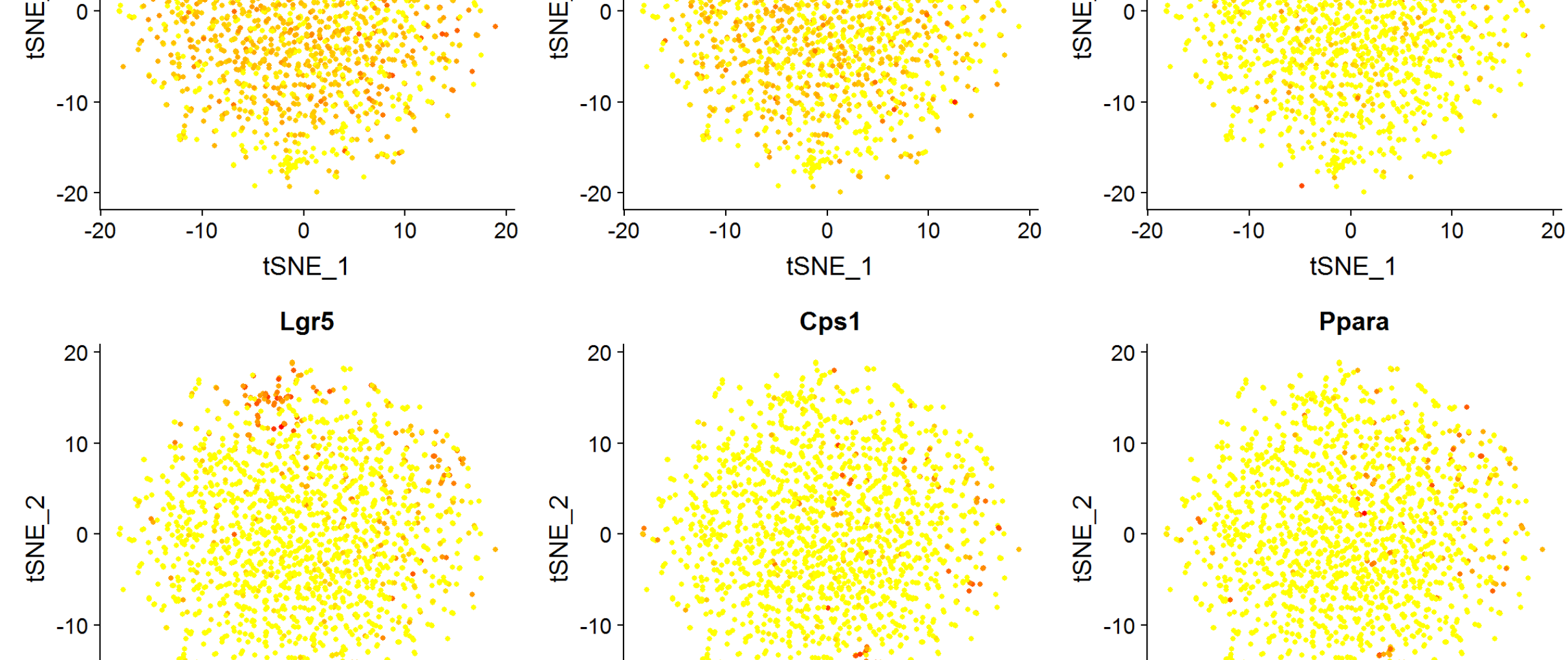


```
cdscomb <- FindClusters(cdscomb, reduction.type = "cca,aligned",
resolution = .7, dims.use = 1:20)
p2 <- TSNEPlot(cdscomb, do.label = T, do.return = T, pt.size = 0.5)
plot_grid(p1, p2)
```



We see that there are 4 clusters found, and we look at markers that identify hepatic expression.

```
cluster<-c("Tbx3", "Id3", "Etv5", "Lgr5"), "Hepatoblast genes"
cluster<-c("Cps1", "Ppara", "ApoA", "Cyp2a26"), "Hepatocyte c genes"
FeaturePlot(object = cdscomb, features.plot=c(cluster,cluster2))
```



We see that cluster 3 can be identified as hepatocytes

Conserved markers are identified and plotted in a similar fashion, we use cluster 3 as an example.

```
nk.markers <- FindConservedMarkers(cdscomb, ident.1=3, grouping.var = "protocol",
print.bar = FALSE)
```

```
## Warning in if (ident.use.2 %in% object$ident) {: the condition has length
## > 1 and only the first element will be used
## Warning in if (ident.use.2 %in% object$ident) {: the condition has length
## > 1 and only the first element will be used
```

```
ccamarkers<-rownames(nk.markers)
FeaturePlot(object = cdscomb, features.plot=ccamarkers[1:12])

#Gm4951 Fgl1 Serpina1e Gc
Hsd3b3 Pzp Cyp2a22 Pemt
Serpina1c Abca8a Itih4 Cyp4a12a
```

```
write.table(ccamarkers, "ccamarkers.txt")
tsnevals<-as.matrix(cdscomb@dr$tne@cell.embeddings)
write.xlsx(tsnevals, "tsnevals.xlsx")
ccavals<-as.matrix(cdscomb@dr$cca@cell.embeddings)
ccavals<-ccavals[,1:2]
write.xlsx(ccavals, "ccavals.xlsx")
ccamarkers<-as.matrix(cdscomb@dr$cca@gene.loadings.full)
ccamarkers<-ccamarkers[,1:2]
write.xlsx(ccamarkers, "ccamarkers.xlsx")
ccamarkers<-as.matrix(cdscomb@dr$cca@aligned@cell.embeddings)
write.xlsx(ccamarkers, "ccamarkers.xlsx")
```