# Hierarchical Clustering

Roger D. Peng, Associate Professor of Biostatistics

May 18, 2016

# Can we find things that are close together?

Clustering organizes things that are **close** into groups

- How do we define close?
- How do we group things?
- How do we visualize the grouping?
- How do we interpret the grouping?

# Hugely important/impactful



```
http://scholar.google.com/scholar?hl=en&q=cluster+
analysis&btnG=&as_sdt=1%2C21&as_sdtp=
```

# Hierarchical clustering

- An agglomerative approach
- Find closest two things
- Put them together
- Find next closest
- Requires
- A defined distance
- A merging approach
- Produces
- A tree showing how close things are to each other

# How do we define close?

- Most important step
- Garbage in -> garbage out
- Distance or similarity
- Continuous - euclidean distance
- Continuous - correlation similarity
- Binary - manhattan distance
- Pick a distance/similarity that makes sense for your problem

# Example distances - Euclidean

# Example distances - Euclidean

$$\sqrt{(X_1 - X_2)^2 + (Y_1 - Y_2)^2}$$

Baltimore $(X_2, Y_2)$

$(Y_1 - Y_2)$

DC
$(X_1, Y_1)$

$(X_1 - X_2)$

In general:

# Example distances - Manhattan



In general:

# Hierarchical clustering - example

```r
set.seed(1234); par(mar=c(0,0,0,0))
x <- rnorm(12,mean=rep(1:3,each=4),sd=0.2)
y <- rnorm(12,mean=rep(c(1,2,1),each=4),sd=0.2)
plot(x,y,col="blue",pch=19,cex=2)
text(x+0.05,y+0.05,labels=as.character(1:12))
```

# Hierarchical clustering - dist

- ▶ Important parameters: *x,method*

```
dataFrame <- data.frame(x=x,y=y)
dist(dataFrame)
```

```
##             1          2          3          4
## 2  0.34120511
## 3  0.57493739 0.24102750
## 4  0.26381786 0.52578819 0.71861759
## 5  1.69424700 1.35818182 1.11952883 1.80666768
## 6  1.65812902 1.31960442 1.08338841 1.78081321 0.0815026
## 7  1.49823399 1.16620981 0.92568723 1.60131659 0.2111043
## 8  1.99149025 1.69093111 1.45648906 2.02849490 0.6170420
## 9  2.13629539 1.83167669 1.67835968 2.35675598 1.1834965
## 10 2.06419586 1.76999236 1.63109790 2.29239480 1.2384787
## 11 2.14702468 1.85183204 1.71074417 2.37461984 1.2815394
## 12 2.05664233 1.74662555 1.58658782 2.27232243 1.0770097
##             7          8          9          10
```
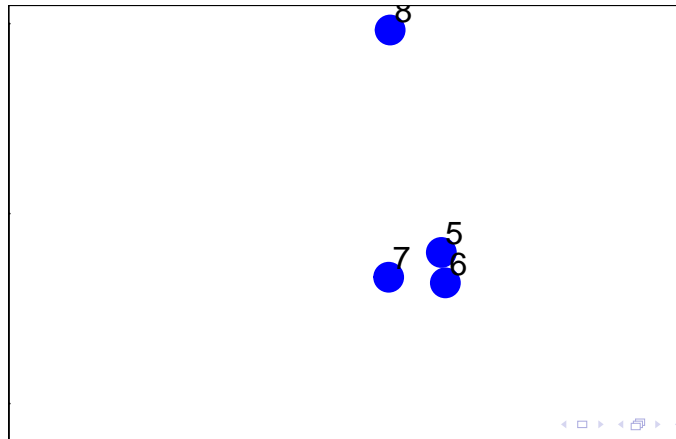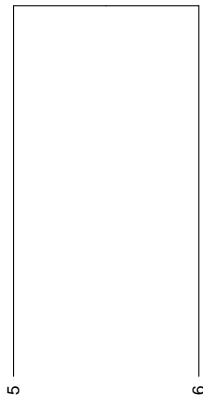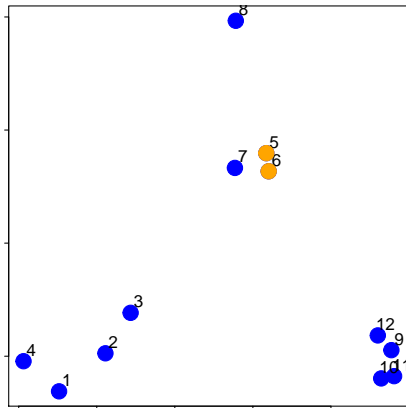
# Hierarchical clustering - #1

# Hierarchical clustering - #2

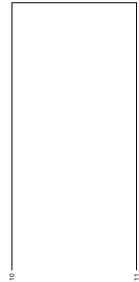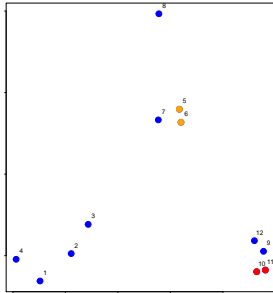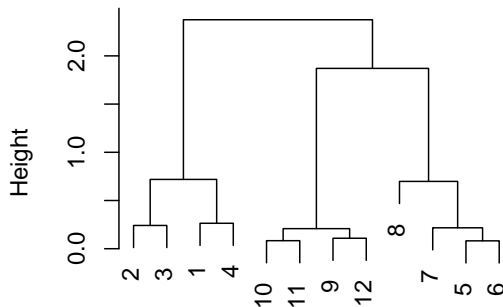# Hierarchical clustering - #3

# Hierarchical clustering - hclust

```
dataFrame <- data.frame(x=x,y=y)
distxy <- dist(dataFrame)
hClustering <- hclust(distxy)
plot(hClustering)
```
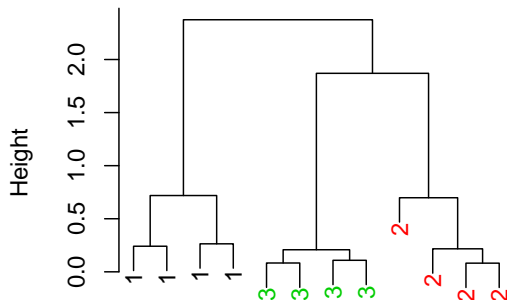
**Cluster Dendrogram**

# Prettier dendrograms

```
myplclust <- function( hclust, lab=hclust$labels, lab.col=
  ## modifiction of plclust for plotting hclust objects *i
  ## Copyright Eva KF Chan 2009
  ## Arguments:
  ##    hclust:    hclust object
  ##    lab:       a character vector of labels of the lea
  ##    lab.col:   colour for the labels; NA=default devic
  ##    hang:      as in hclust & plclust
  ## Side effect:
  ##    A display of hierarchical cluster with coloured lea
  y <- rep(hclust$height,2); x <- as.numeric(hclust$merge)
  y <- y[which(x<0)]; x <- x[which(x<0)]; x <- abs(x)
  y <- y[order(x)]; x <- x[order(x)]
  plot( hclust, labels=FALSE, hang=hang, ... )
  text( x=x, y=y[hclust$order]-(max(hclust$height)*hang),
        labels=lab[hclust$order], col=lab.col[hclust$order]
        srt=90, adj=c(1,0.5), xpd=NA, ... )
}
```
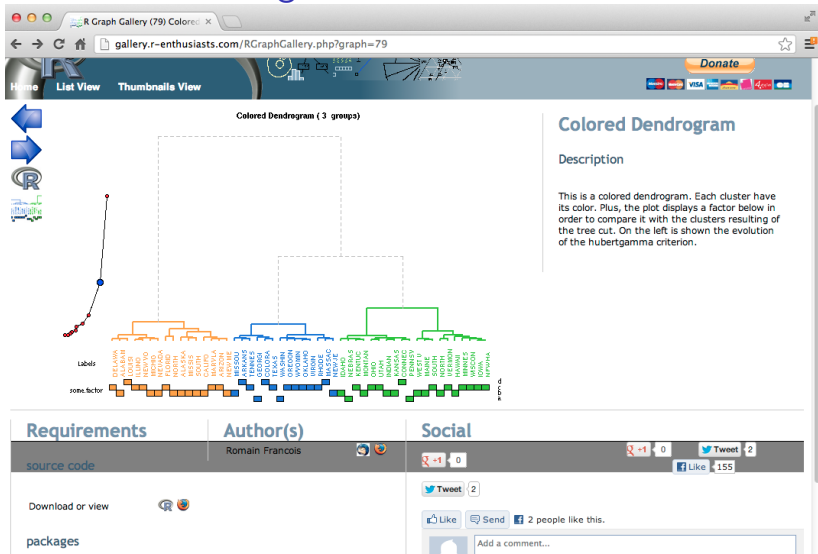
# Pretty dendrograms

```
dataFrame <- data.frame(x=x,y=y)
distxy <- dist(dataFrame)
hClustering <- hclust(distxy)
myplclust(hClustering,lab=rep(1:3,each=4),lab.col=rep(1:3,e
```
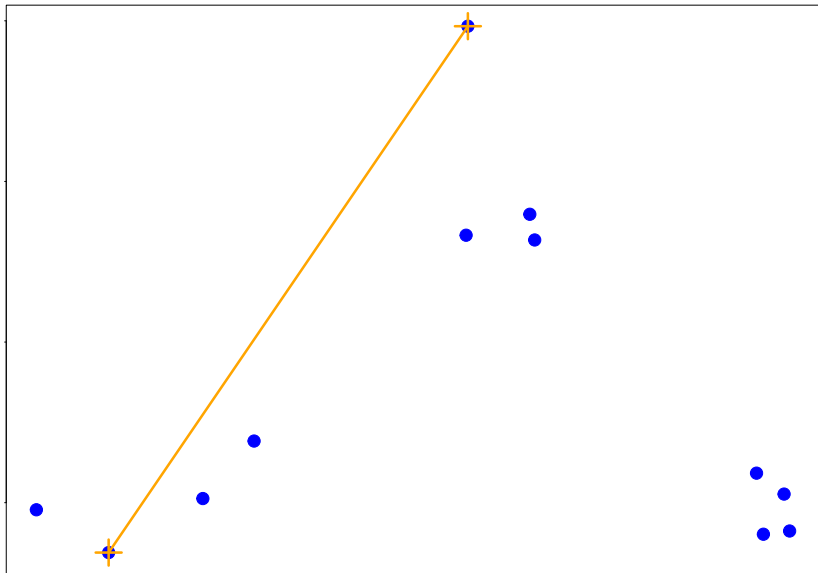


**Cluster Dendrogram**

# Even Prettier dendrograms



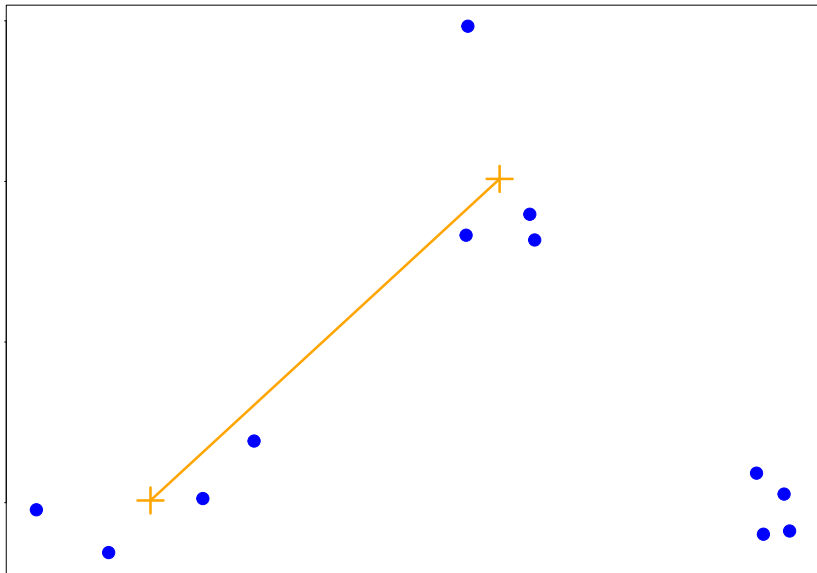http://gallery.r-enthusiasts.com/RGraphGallery.php?
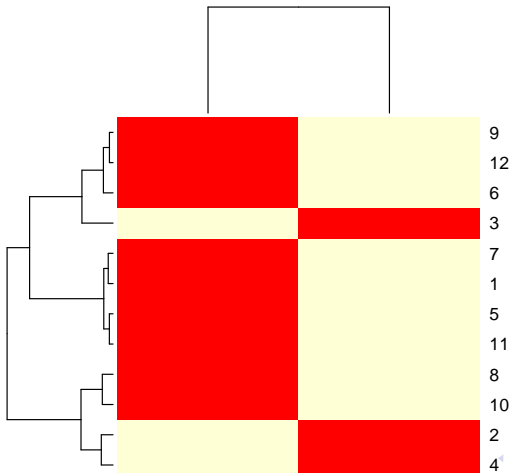graph=79

# Merging points - complete

# Merging points - average

# heatmap()

```
dataFrame <- data.frame(x=x,y=y)
set.seed(143)
dataMatrix <- as.matrix(dataFrame)[sample(1:12),]
heatmap(dataMatrix)
```

# Notes and further resources

- Gives an idea of the relationships between variables/observations
- The picture may be unstable
- Change a few points
- Have different missing values
- Pick a different distance
- Change the merging strategy
- Change the scale of points for one variable
- But it is deterministic
- Choosing where to cut isn't always obvious
- Should be primarily used for exploration
- Rafa's Distances and Clustering Video
- Elements of statistical learning