

# ISSR Short Course

Gregory J. Matthews <sup>1</sup>

<sup>1</sup>Department of Mathematics and Statistics  
Loyola University Chicago

June 2015

# Outline

Basic Descriptives

Principal Component Analysis

Cluster Analysis

# Jaccard Index

Similarity

$$J(A, B) = \frac{|A \cap B|}{|A \cup B|}$$

Dissimilarity

$$J(A, B) = \frac{|A \cup B| - |A \cap B|}{|A \cup B|}$$

```
library(tm)

## Loading required package: NLP

#Example of jaccard index
example<-list()
example[[1]]<-PlainTextDocument("The cat dog dog")
example[[2]]<-PlainTextDocument("The dog")
exampleCorp<-Corpus(VectorSource(example))
exampleTDM<-TermDocumentMatrix(exampleCorp)
proxy::dist(as.matrix(t(exampleTDM)), method = "Jaccard")

##           1
## 2 0.3333333
```

```
x<-c("The","cat","dog","dog")
y<-c("The","dog")
#Jaccard Index
1-length(intersect(x,y))/length(union(x,y))

## [1] 0.3333333
```

```
library(tm)
#Example of jaccard index
example<-list()
example[[1]]<-PlainTextDocument("The cat dog dog")
example[[2]]<-PlainTextDocument("The dog")
exampleCorp<-Corpus(VectorSource(example))
exampleTDM<-TermDocumentMatrix(exampleCorp)
proxy::dist(as.matrix(t(exampleTDM)), method = "eJaccard")

##      1
## 2 0.4
```

```
x<-c(1,2,1)
y<-c(0,1,1)
#eJaccard Index
1-x%%y/(x%%x+y%%y-x%%y)

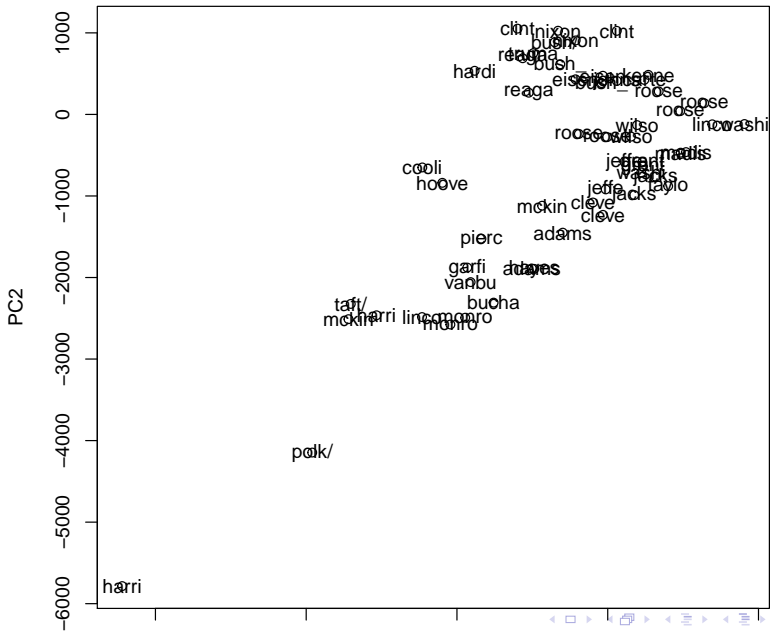
##          [,1]
## [1,]      0.4
```

- ▶ Principal Component Analysis (PCA) is a dimensionality reduction technique.



```
#Principal Component Analysis
corr<-cor(as.matrix(presTDM))
evv<-eigen(corr)
pcs <- evv$eigenvectors[,1:2]
##values:a vector containing the p eigenvalues of x, sorted in decreasing order
evals <- evv$values[1:2]
temp<-diag(1/sqrt(evals)) %*% t(pcs)%*%t(as.matrix(presTDM))
PCs<-t(as.matrix(presTDM))%*%t(temp)
```

## PC1 vs PC2: Pres Inaugural Speeches



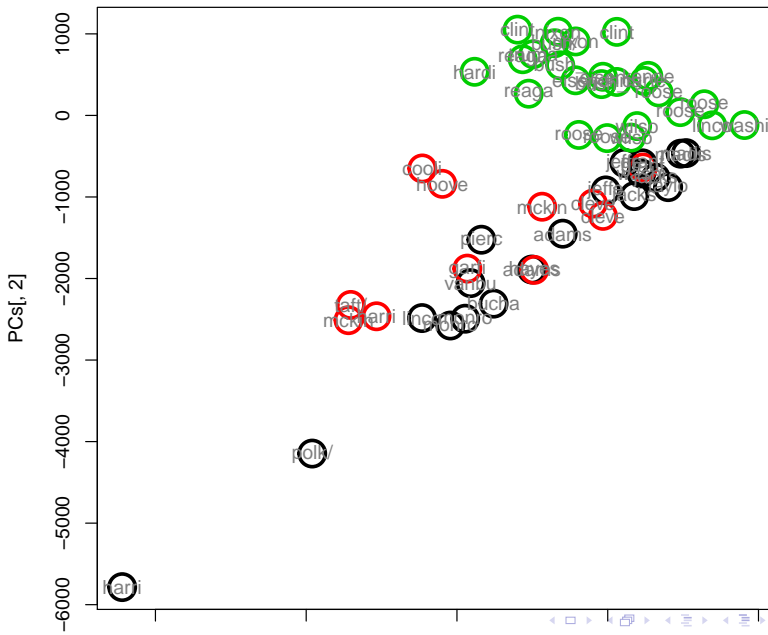
- ▶ Cluster analysis allows us to cluster objects into one of  $K$  clusters where  $K$  can be either known or unknown.
- ▶ We'll look at a simple clustering method called  $k$ -means.
- ▶ Essentially, clusters are created and the center of the cluster is calculated.
- ▶ Objects end up in the cluster so that the distances between the object and the cluster centroid are minimized.

- ▶ Let's go back to the presidential inaugural speeches and cluster those.
- ▶ We will start with 2 clusters and then explore the possibility of more clusters.

```
#Term frequency - Inverse document frequency
m<-weightTfIdf(presTDM)
d<-proxy::dist(as.matrix(t(m)), method = "eJaccard")
cl <- kmeans(d, 3)
table(cl$cluster)

##
##  1  2  3
## 19 11 25
```

## PC1 vs PC2: k-means clustering



```
m<-weightTfIdf(presTDM)
d<-proxy::dist(as.matrix(t(m)), method = "eJaccard")
hc <- hclust(d, method="average")
```

# Cluster Dendrogram

