# Statistics 452: Statistical Learning and Prediction

## Chapter 10: Introduction to Unsupervised Learning

Brad McNeney

2017-11-16

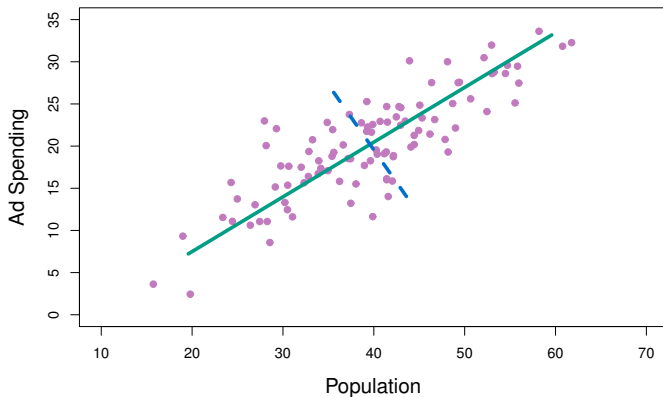# Supervised *versus* Unsupervised Learning

- ▶ Supervised means that there is an outcome **y**, unsupervised means there is not.
- ▶ Supervised learning has well-defined goals like prediction.
    - ▶ Can check the fitted model by seeing how well it predicts test observations.
- ▶ Unsupervised learning is more exploratory, without an obvious goal.
    - ▶ A common theme is trying to identify simple structure underlying the feature data.
    - ▶ We will discuss dimension reduction by principal components analysis (PCA) and clustering.

# Principal Components Analysis (PCA)

- ▶ Goal is low-rank approximation of the $X$ data matrix
  - ▶ Discussed in Chapter 6 and reviewed below.
- ▶ Think of principal components (PCs) as new coordinates for the data vectors.
  - ▶ The first PC is the direction of greatest variation,
  - ▶ The second PC is the direction of second-greatest variation, orthogonal to the first,
  - ▶ And so on.

# PCs for Advertising Data

- Text Figure 6.14: The green line is the first PC, the blue line the second.

# PCs as Linear Combinations of $X$'s

- The details of how the linear combinations are derived are discussed in the text.
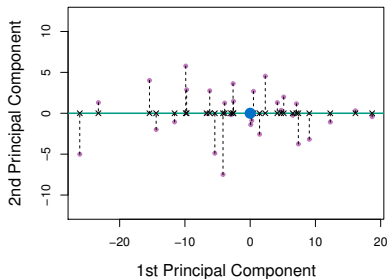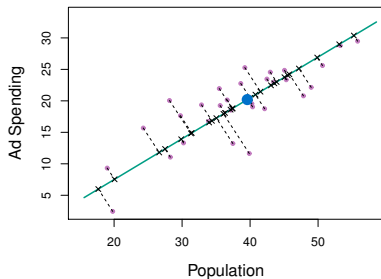- In the advertising example, the first PC is

$$Z_1 = 0.838X_1 + 0.544X_2$$

where $X_1$ is population centred by its mean and $X_2$ is advertising expenditure centred by its mean.

- The coefficients of the linear combination, $\phi_{11} = 0.838$ and $\phi_{12} = 0.544$, are called the first principal component *loadings*.
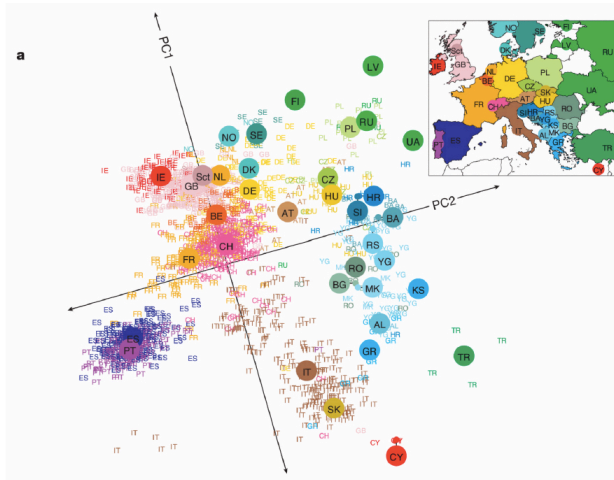
# Principal Component Scores

- ▶ Projecting each point onto the PCs gives the PC scores.
  - ▶ Projecting a data vector onto a line means finding the point on the line closest to the vector.

- ▶ Text Figure 6.15: Black x's are the first PC score for each observation, distance of each purple dot from the green line is the second PC score.

# High-Dimensional Example: Genes Reflect Geography

- First 2 PCs from 197,146 genetic markers on 1,387 European individuals (Novembre *et al.* 2008)

# US Arrests Data

- ▶ Dataset that comes with R.
- ▶ From the help file: "This data set contains statistics, in arrests per 100,000 residents for assault, murder, and rape in each of the 50 US states in 1973. Also given is the percent of the population living in urban areas."

```
data(USArrests) # help(USArrests)
head(USArrests)
```

```
##            Murder Assault UrbanPop Rape
## Alabama      13.2     236       58 21.2
## Alaska       10.0     263       48 44.5
## Arizona       8.1     294       80 31.0
## Arkansas      8.8     190       50 19.5
## California    9.0     276       91 40.6
## Colorado      7.9     204       78 38.7
```

```r
pcout <- prcomp(USArrests,scale=TRUE)
pcout$rotation # loadings
```

```
##                PC1        PC2        PC3         PC4
## Murder   -0.5358995  0.4181809 -0.3412327  0.64922780
## Assault  -0.5831836  0.1879856 -0.2681484 -0.74340748
## UrbanPop -0.2781909 -0.8728062 -0.3780158  0.13387773
## Rape     -0.5434321 -0.1673186  0.8177779  0.08902432
```
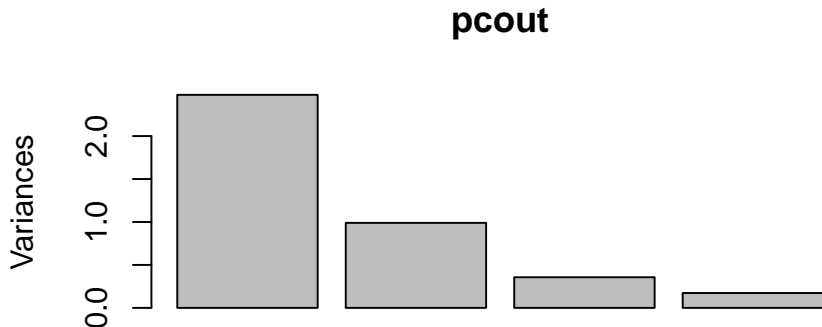
```r
head(pcout$x) # scores
```

```
##                   PC1        PC2         PC3          PC4
## Alabama    -0.9756604  1.1220012 -0.43980366  0.154696581
## Alaska     -1.9305379  1.0624269  2.01950027 -0.434175454
## Arizona    -1.7454429 -0.7384595  0.05423025 -0.826264240
## Arkansas    0.1399989  1.1085423  0.11342217 -0.180973554
## California -2.4986128 -1.5274267  0.59254100 -0.338559240
## Colorado   -1.4993407 -0.9776297  1.08400162  0.001450164
```

# Scree Plot

- A scree plot shows the variance (or proportion of total variance) in the direction of each PC.
- If the variance drops and then levels out, the "elbow" where it levels out is a reasonable choice for a reduced number of PCs that captures most of the variation in the **X**.

```
screeplot(pcout) # or just plot(pcout)
```



**pcout**

- No obvious elbow.

# Iris Data
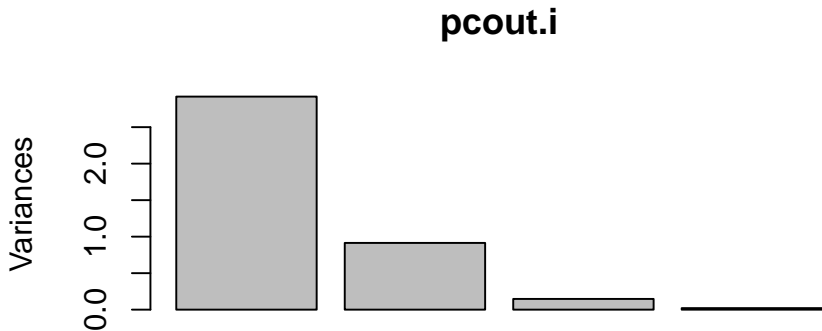
```
data(iris) # help(iris)
head(iris)
```

```
##   Sepal.Length Sepal.Width Petal.Length Petal.Width Species
## 1          5.1         3.5          1.4         0.2  setosa
## 2          4.9         3.0          1.4         0.2  setosa
## 3          4.7         3.2          1.3         0.2  setosa
## 4          4.6         3.1          1.5         0.2  setosa
## 5          5.0         3.6          1.4         0.2  setosa
## 6          5.4         3.9          1.7         0.4  setosa
```

```
pcout.i <- prcomp(iris[,-5],scale=TRUE)
```

- For the iris data, two PCs appear to explain most of the variation.

```
screeplot(pcout.i)
```



**pcout.i**

# Interpretation of Loadings

- The first PC is a contrast between sepal width and the other variables.
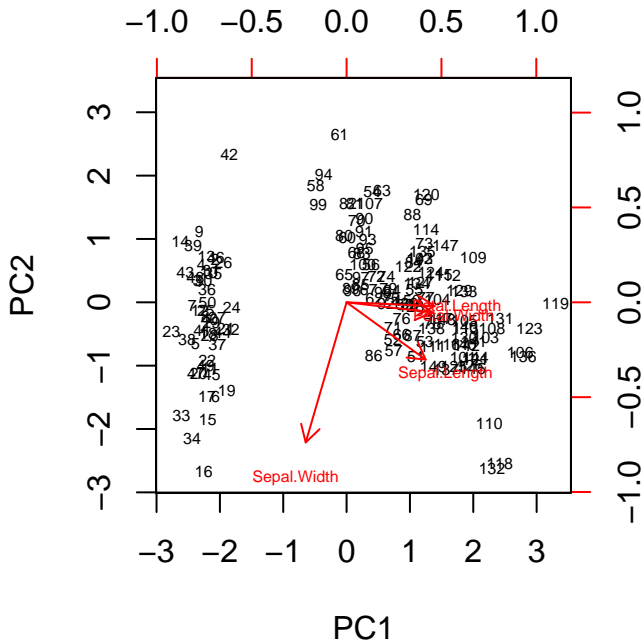- The second PC is a weighted average of sepal length and width.

```
pcout.i$rotation
```

```
##                     PC1         PC2        PC3        PC4
## Sepal.Length  0.5210659 -0.37741762  0.7195664  0.2612863
## Sepal.Width  -0.2693474 -0.92329566 -0.2443818 -0.1235096
## Petal.Length  0.5804131 -0.02449161 -0.1421264 -0.8014492
## Petal.Width   0.5648565 -0.06694199 -0.6342727  0.5235971
```
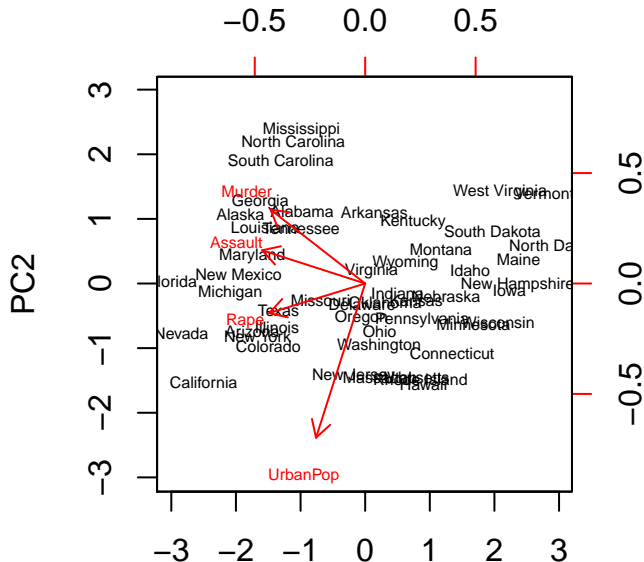
# Biplot of First Two PCs

- ▶ We can visualize the first two PCs on a scatterplot.
- ▶ A biplot shows
    1. the PC scores for observational units, and
    2. the loadings of the features that define the first two PCs

```
biplot(pcout.i,cex=.5,scale=0) #scale=0 avoids scaling of points on plot
```

# Biplot of US Arrests Data

```
biplot(pcout,cex=.5,scale=0) #scale=0 avoids scaling of points on plot
```

- Note different appearance from text: PCs are only unique up to a sign change