# Lecture 3: Principal Components Analysis (PCA)

### Reading: Sections 6.3.1, 10.1, 10.2, 10.4

**STATS 202: Data mining and analysis**

Sergio Bacallado
October 7, 2013

# Announcements

- Kaggle invitations have been sent. You have to register + join the competition.

- If you want to form a team, you should do so before making any submissions.

- I'm extending my office hours on Monday for Python help.

- Warning: Some nodes on corn have a different version of IPython.

# The bias variance decomposition

Let $x_0$ be fixed, $y_0 = f(x_0) + \varepsilon$, and $\hat{f}$ be estimated from $n$ separate training samples $\{x_1, \ldots, x_n; y_1, \ldots, y_n\}$.

Let $E$ denote the expectation over $\varepsilon$ and the training samples. Then, the Mean Squared Error at $x_0$ can be decomposed:

$$MSE(x_0) = E(y_0 - \hat{f}(x_0))^2 = \mathsf{Var}(\hat{f}(x_0)) + [\mathsf{Bias}(\hat{f}(x_0))]^2 + \mathsf{Var}(\varepsilon).$$

# The bias variance decomposition

Let $x_0$ be fixed, $y_0 = f(x_0) + \varepsilon$, and $\hat{f}$ be estimated from $n$ separate training samples $\{x_1, \ldots, x_n; y_1, \ldots, y_n\}$.

Let $E$ denote the expectation over $\varepsilon$ and the training samples. Then, the Mean Squared Error at $x_0$ can be decomposed:

$$MSE(x_0) = E(y_0 - \hat{f}(x_0))^2 = \mathsf{Var}(\hat{f}(x_0)) + [\mathsf{Bias}(\hat{f}(x_0))]^2 + \mathsf{Var}(\varepsilon).$$

Irreducible error

# The bias variance decomposition

Let $x_0$ be fixed, $y_0 = f(x_0) + \varepsilon$, and $\hat{f}$ be estimated from $n$ separate training samples $\{x_1, \ldots, x_n; y_1, \ldots, y_n\}$.

Let $E$ denote the expectation over $\varepsilon$ and the training samples. Then, the Mean Squared Error at $x_0$ can be decomposed:

$$MSE(x_0) = E(y_0 - \hat{f}(x_0))^2 = \mathsf{Var}(\hat{f}(x_0)) + [\mathsf{Bias}(\hat{f}(x_0))]^2 + \mathsf{Var}(\varepsilon).$$

The variance of the estimate of $Y$: $E[\hat{f}(x_0) - E(\hat{f}(x_0))]^2$

This measures how much the estimate of $\hat{f}$ at $x_0$ changes when we sample new training data.

# The bias variance decomposition

Let $x_0$ be fixed, $y_0 = f(x_0) + \varepsilon$, and $\hat{f}$ be estimated from $n$ separate training samples $\{x_1, \ldots, x_n; y_1, \ldots, y_n\}$.

Let $E$ denote the expectation over $\varepsilon$ and the training samples. Then, the Mean Squared Error at $x_0$ can be decomposed:

$$MSE(x_0) = E(y_0 - \hat{f}(x_0))^2 = \mathsf{Var}(\hat{f}(x_0)) + [\mathsf{Bias}(\hat{f}(x_0))]^2 + \mathsf{Var}(\varepsilon).$$

The squared bias of the estimate of $Y$: $[E(\hat{f}(x_0) - f(x_0))]^2$

This measures the deviation of the average estimate $\hat{f}$ at $x_0$ from $f(x_0)$.

# The bias variance decomposition

Let $x_0$ be fixed, $y_0 = f(x_0) + \varepsilon$, and $\hat{f}$ be estimated from $n$ separate training samples $\{x_1, \ldots, x_n; y_1, \ldots, y_n\}$.

Let $E$ denote the expectation over $\varepsilon$ and the training samples. Then, the Mean Squared Error at $x_0$ can be decomposed:

$$MSE(x_0) = E(y_0 - \hat{f}(x_0))^2 = \mathsf{Var}(\hat{f}(x_0)) + [\mathsf{Bias}(\hat{f}(x_0))]^2 + \mathsf{Var}(\varepsilon).$$

Both variance and squared bias are always positive.

Higher variance $\iff$ More flexibility $\iff$ Lower bias.

We will aim to minimize both sources of error simultaneously.
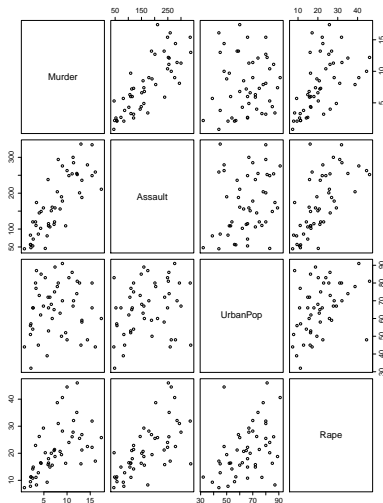
# Principal Components Analysis

- This is the most popular unsupervised procedure ever.

- Invented by Karl Pearson (1901).

- Developed by Harold Hotelling (1933).

- **What does it do?** It provides a way to visualize high dimensional data, summarizing the most important information.

# Principal Components Analysis

- This is the most popular unsupervised procedure ever.

- Invented by Karl Pearson (1901).

- Developed by Harold Hotelling (1933). ← Stanford pride!

- **What does it do?** It provides a way to visualize high dimensional data, summarizing the most important information.
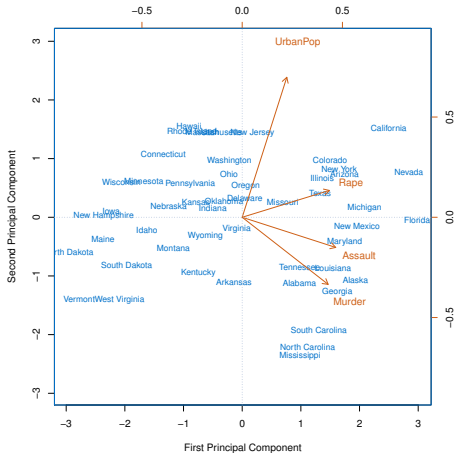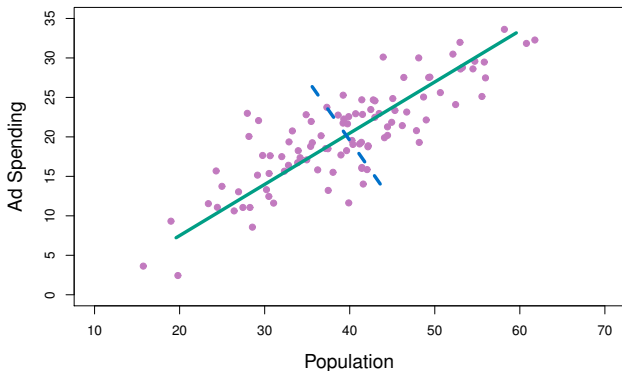
# What is PCA good for?

# What is PCA good for?



Figure 10.1

# What is the first principal component?

It is the vector which passes the closest to a cloud of samples, in terms of Euclidean distance.

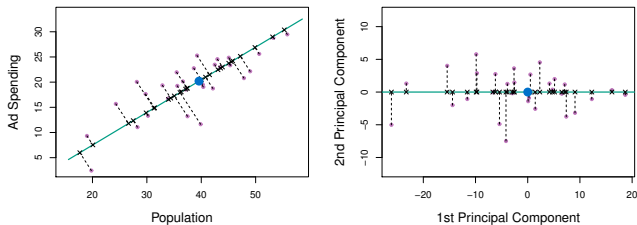i.e. The green direction minimizes the average length of the dotted lines.



Figure 6.15

# What does this look like with 3 variables?

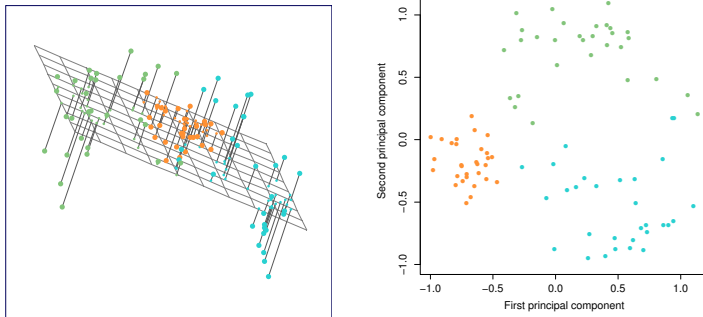The first two principal components span a plane which is closest to the data.



Figure 10.2

# A second interpretation

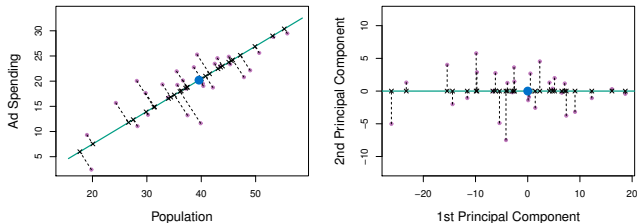The projection onto the first principal component is the one with the **highest variance**.



Figure 6.15

# How do we say this in math?

Let $\mathbf{X}$ be a data matrix with $n$ samples, and $p$ variables. From each variable, we subtract the mean of the column; i.e. we **center** the variables.

To find the first principal component $\phi_1 = (\phi_{11}, \ldots, \phi_{p1})$, we solve the following optimization

$$\max_{\phi_{11}, \ldots, \phi_{p1}} \left\{ \frac{1}{n} \sum_{i=1}^{n} \left( \sum_{j=1}^{p} \phi_{j1} x_{ij} \right)^2 \right\}$$

$$\text{subject to} \sum_{j=1}^{p} \phi_{j1}^2 = 1.$$

# How do we say this in math?

Let $\mathbf{X}$ be a data matrix with $n$ samples, and $p$ variables. From each variable, we subtract the mean of the column; i.e. we **center** the variables.

To find the first principal component $\phi_1 = (\phi_{11}, \ldots, \phi_{p1})$, we solve the following optimization

$$\max_{\phi_{11},\ldots,\phi_{p1}} \left\{ \frac{1}{n} \sum_{i=1}^{n} \left( \sum_{j=1}^{p} \phi_{j1} x_{ij} \right)^2 \right\}$$

$$\text{subject to } \sum_{j=1}^{p} \phi_{j1}^2 = 1.$$

Projection of the $i$th sample onto $\phi_1$. Also known as **the score** $z_{i1}$

# How do we say this in math?

Let $\mathbf{X}$ be a data matrix with $n$ samples, and $p$ variables. From each variable, we subtract the mean of the column; i.e. we **center** the variables.

To find the first principal component $\phi_1 = (\phi_{11}, \ldots, \phi_{p1})$, we solve the following optimization

$$\max_{\phi_{11},\ldots,\phi_{p1}} \left\{ \frac{1}{n} \sum_{i=1}^{n} \left( \sum_{j=1}^{p} \phi_{j1} x_{ij} \right)^2 \right\}$$

$$\text{subject to } \sum_{j=1}^{p} \phi_{j1}^2 = 1.$$

Variance of the $n$ samples projected onto $\phi_1$.

# How do we say this in math?

To find the second principal component $\phi_2 = (\phi_{12}, \ldots, \phi_{p2})$, we solve the following optimization

$$\max_{\phi_{12}, \ldots, \phi_{p2}} \left\{ \frac{1}{n} \sum_{i=1}^{n} \left( \sum_{j=1}^{p} \phi_{j2} x_{ij} \right)^2 \right\}$$

subject to $\sum_{j=1}^{p} \phi_{j2}^2 = 1$ and $\sum_{j=1}^{p} \phi_{j1}\phi_{j2} = 0$.

First and second principal components must be orthogonal.

Equivalent to saying that the scores $(z_{11}, \ldots, z_{n1})$ and $(z_{12}, \ldots, z_{n2})$ are uncorrelated.

# Solving the optimization

This optimization is fundamental in linear algebra. It is satisfied by either:

- The singular value decomposition (SVD) of $\mathbf{X}$:

$$\mathbf{X} = \mathbf{U}\mathbf{\Sigma}\mathbf{\Phi}^T$$

  where the $i$th column of $\mathbf{\Phi}$ is the $i$th principal component $\phi_i$, and the $i$th column of $\mathbf{U}\mathbf{\Sigma}$ is the $i$th vector of scores $(z_{1i}, \ldots, z_{ni})$.

- The eigendecomposition of $\mathbf{X}^T\mathbf{X}$:

$$\mathbf{X}^T\mathbf{X} = \mathbf{\Phi}\mathbf{\Sigma}^2\mathbf{\Phi}^T$$
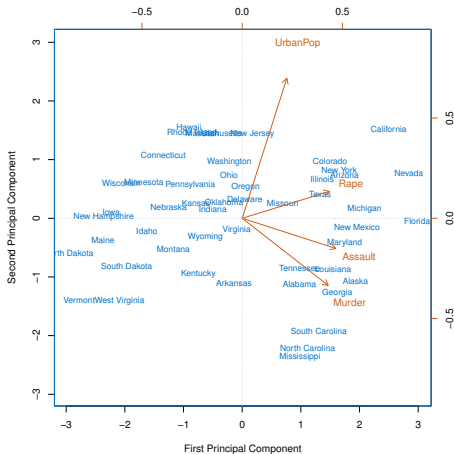
# PCA in practice: The biplot



Figure 10.1

# Scaling the variables

Most of the time, we don't care about the absolute numerical value of a variable. We care about the value relative to the spread observed in the sample.

Before PCA, in addition to **centering** each variable, we also multiply it times a constant to make its variance equal to 1.
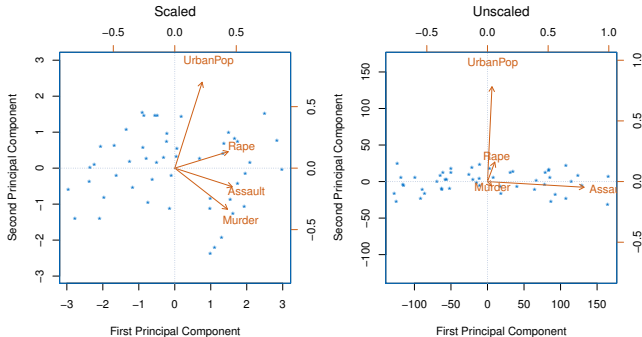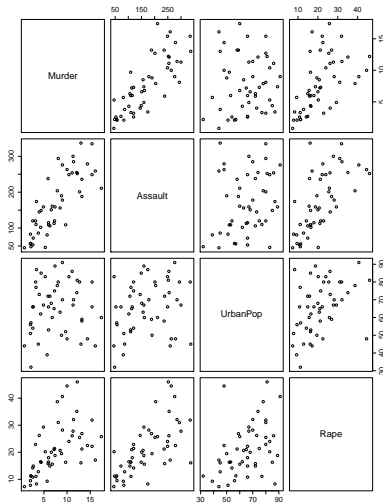
# Example: scaled vs. unscaled PCA



Figure 10.3

# Scaling the variables

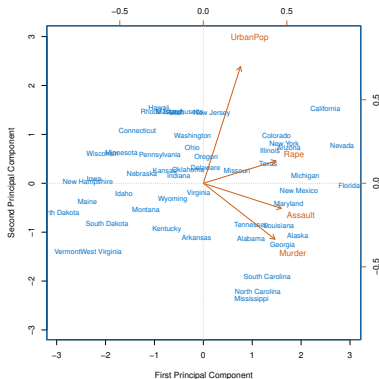In special cases, we have variables measured in the same unit; e.g. gene expression levels for different genes.

Therefore, we care about the absolute value of the variables and we can perform PCA without scaling.

# How many principal components are enough?

# How many principal components are enough?



We said 2 principal components capture most of the relevant information. But how can we tell?

# The proportion of variance explained

We can think of the top **principal components** as directions in space in which the data vary the most.

The $i$th **score vector** $(z_{1i}, \ldots, z_{ni})$ can be interpreted as a *new* variable. The variance of this variable decreases as we take $i$ from 1 to $p$. However, the total variance of the score vectors is the same as the total variance of the original variables:
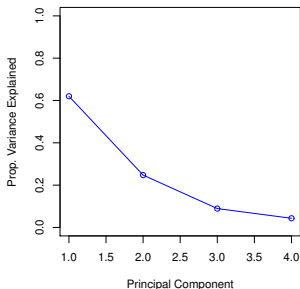
$$\sum_{k=1}^{p} \frac{1}{n} \sum_{j=1}^{n} x_{jk}^2.$$

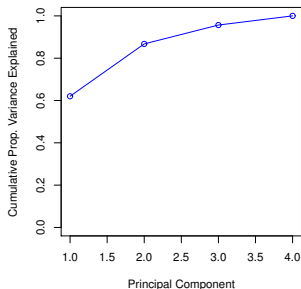We can quantify how much of the variance is captured by the first $m$ principal components/score variables.

# The proportion of variance explained

The variance of the $m$th score variable is:

$$\frac{1}{n}\sum_{i=1}^{n} z_{im}^2 = \frac{1}{n}\sum_{i=1}^{n}\left(\sum_{j=1}^{p}\phi_{jm}x_{ij}\right)^2 = \frac{1}{n}\mathbf{\Sigma}_{mm}^2.$$



Scree plot

# Generalizations of PCA

PCA works under a Euclidean geometry in the space of variables. Often, the natural geometry is different:

- ▶ We expect some variables to be "closer" to each other that to other variables.
- ▶ Some correlations between variables would be more surprising than others.

Examples:

- ▶ Variables are pixel values, samples are different images of the brain. We expect neighboring pixels to have stronger correlations.
- ▶ Variables are counts for different gut bacteria, samples are different people. We expect bacteria that are close phylogenetically to be strongly correlated.

# Generalizations of PCA

There are ways to include this knowledge in a PCA. See:

1. Susan Holmes. *Multivariate Analysis, the French way.* (2006).
2. Omar de la Cruz and Susan Holmes. *An introduction to the duality diagram.* (2011).
3. Stéphane Dray and Thibaut Jombart. *Revisiting GuerryŐs data: Introducing spatial constraints in multivariate analysis.* (2011).
4. Genevera Allen, Logan Grosenick, and Jonathan Taylor. *A Generalized Least Squares Matrix Decomposition.* (2011).