

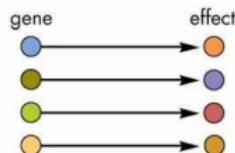
Latent Variable Models

Schwartz

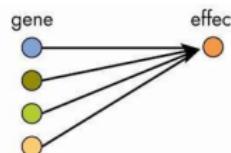
August 13, 2017

Genetics Primer

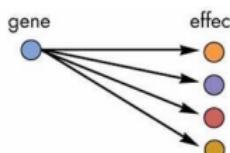
- Polygenic Inheritance: trait is controlled by 2+ genes
- Pleiotropy: gene drives 2+ seemingly unrelated traits



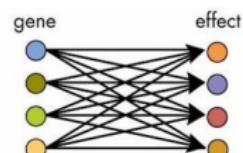
Each gene has a distinct
biological effect.



Polygenic trait: Many genes
contribute to a single effect.



Pleiotropy: A gene has
multiple effects.



Polygenic traits and pleiotropy

$$\begin{bmatrix} \sigma_1^2 & 0 & \cdots & 0 \\ 0 & \sigma_2^2 & \cdots & 0 \\ \vdots & \vdots & \ddots & \vdots \\ 0 & 0 & \cdots & \sigma_p^2 \end{bmatrix}$$

Experimental Design

$$f(Y|\mathbf{X})$$

Predictive
Modeling

$$\begin{bmatrix} 1^2 & r_{12} & \cdots & r_{1p} \\ r_{21} & 1^2 & \cdots & r_{2p} \\ \vdots & \vdots & \ddots & \vdots \\ r_{p1} & r_{p2} & \cdots & 1^2 \end{bmatrix}$$

Correlated Features

PCA
FA
ICA

Objectives

- ▶ Covariance Matrices Vs. Correlation Matrices
 - ▶ And what they do
- ▶ Eigen-vectors/values of Correlation & Covariance Matrices
 - ▶ And what they do:
 - ▶ Proportion of Variance Explained
 - ▶ Skree Plot
 - ▶ Directions
 - ▶ Projecting onto eigenvectors
- ▶ Principal Components Analysis (PCA)
 - ▶ Interpretation
 - ▶ Biplots
 - ▶ Dimensionality Reduction
 - ▶ Visualization
- ▶ Singular Value Decomposition (SVD)
 - ▶ Same, but better than PCA?
- ▶ Factor Analysis (FA)
- ▶ Independent Component Analysis (ICA)

Sample Covariance Matrix

$X_{n \times p}$:

Feature	X_1	X_2	\dots	X_p
n=1				
2				
.				
n				

Sample Covariance Matrix

$X_{n \times p}$:

Feature	X_1	X_2	\dots	X_p
n=1				
2				
\vdots				
n				

the sample *covariance matrix* and *correlation matrix* capture observed feature covariances and correlations, respectively

$\Sigma_{p \times p}$	X_1	X_2	\dots	X_p
X_1	σ_{11}^2	σ_{12}^2		σ_{1p}^2
X_2	σ_{21}^2	σ_{22}^2		σ_{2p}^2
\vdots		\ddots		
X_p	σ_{p1}^2	σ_{p2}^2		σ_{pp}^2

$r_{p \times p}$	X_1	X_2	\dots	X_p
X_1	1	r_{12}		r_{1p}
X_2	r_{21}	1		r_{2p}
\vdots		\ddots		
X_p	r_{p1}	r_{p2}		1

Sample Covariance Matrix

$X_{n \times p}$:

Feature	X_1	X_2	\dots	X_p
n=1				
2				
\vdots				
n				

the sample *covariance matrix* and *correlation matrix* capture observed feature covariances and correlations, respectively

$\Sigma_{p \times p}$	X_1	X_2	\dots	X_p
X_1	σ_{11}^2	σ_{12}^2		σ_{1p}^2
X_2	σ_{21}^2	σ_{22}^2		σ_{2p}^2
\vdots		\ddots		
X_p	σ_{p1}^2	σ_{p2}^2		σ_{pp}^2

$r_{p \times p}$	X_1	X_2	\dots	X_p
X_1	1	r_{12}		r_{1p}
X_2	r_{21}	1		r_{2p}
\vdots		\ddots		
X_p	r_{p1}	r_{p2}		1

For a **centered** samples by features data matrix ($X_{ij} = X_{ij}^* - \bar{X}_j^*$)

$$S = \frac{1}{n-1} \mathbf{X}^T \mathbf{X} = \frac{1}{n-1} \sum_{i=1}^n \mathbf{X}_{i \cdot} \mathbf{X}_{i \cdot}^T$$

$$\text{Cov}(X_j, X_{j'}) = \hat{S}_{jj'} = \frac{1}{n-1} \sum_{i=1}^n X_{ij} X_{ij'}$$

$$\hat{\mathbf{r}} = \hat{\mathbf{D}}_S^{-\frac{1}{2}} \hat{\mathbf{S}} \hat{\mathbf{D}}_S^{-\frac{1}{2}}$$

D_S is the $p \times p$ matrix with same diagonal as S

Sample Covariance Matrix

$\mathbf{X}_{n \times p}$:

Feature	X_1	X_2	\dots	X_p
n=1				
2				
\vdots				
n				

the sample *covariance matrix* and *correlation matrix* capture observed feature covariances and correlations, respectively

$\mathbf{S}_{p \times p}$	X_1	X_2	\dots	X_p
X_1	$\hat{\sigma}_{11}^2$	$\hat{\sigma}_{12}^2$		$\hat{\sigma}_{1p}^2$
X_2	$\hat{\sigma}_{21}^2$	$\hat{\sigma}_{22}^2$		$\hat{\sigma}_{2p}^2$
\vdots		\ddots		
X_p	$\hat{\sigma}_{p1}^2$	$\hat{\sigma}_{p2}^2$		$\hat{\sigma}_{pp}^2$

$\hat{\mathbf{r}}_{p \times p}$	X_1	X_2	\dots	X_p
X_1	1	\hat{r}_{12}		\hat{r}_{1p}
X_2	\hat{r}_{21}	1		\hat{r}_{2p}
\vdots		\ddots		
X_p	\hat{r}_{p1}	\hat{r}_{p2}		1

For a **centered** samples by features data matrix ($X_{ij} = X_{ij}^* - \bar{X}_j^*$)

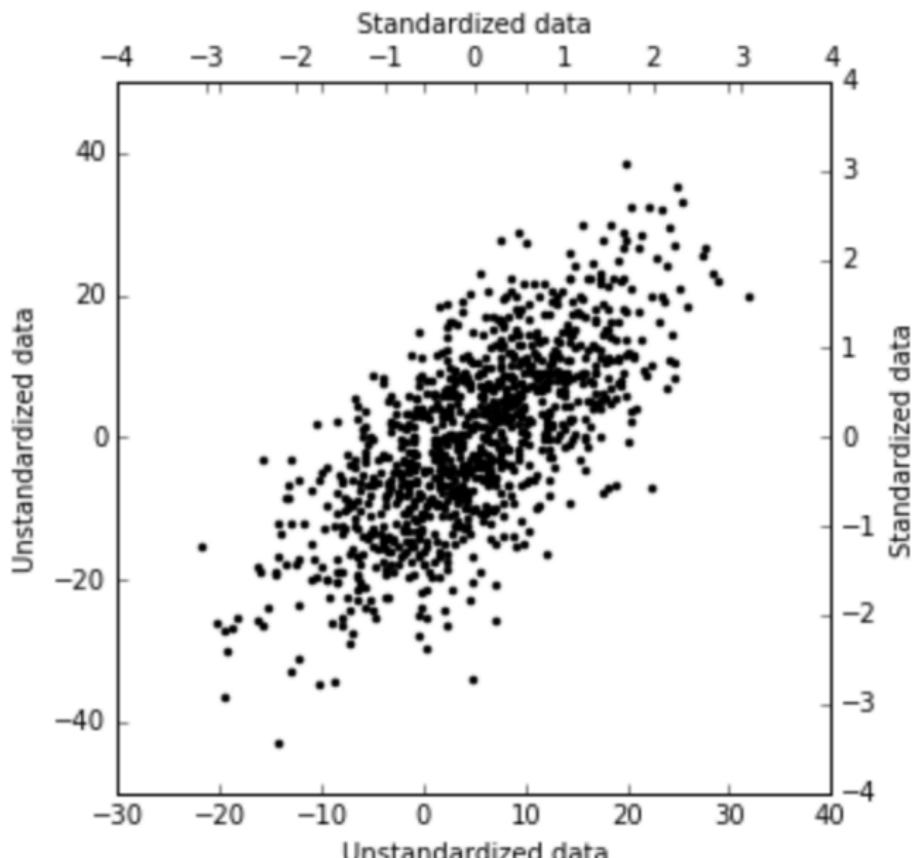
$$\mathbf{S} = \frac{1}{n-1} \mathbf{X}^T \mathbf{X} = \frac{1}{n-1} \sum_{i=1}^n \mathbf{X}_i \mathbf{X}_i^T$$

$$\hat{\mathbf{r}} = \hat{\mathbf{D}}_S^{-\frac{1}{2}} \hat{\mathbf{S}} \hat{\mathbf{D}}_S^{-\frac{1}{2}}$$

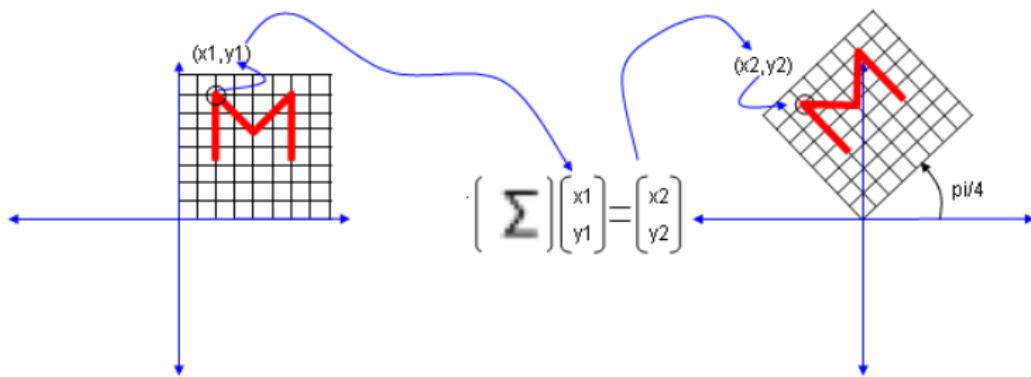
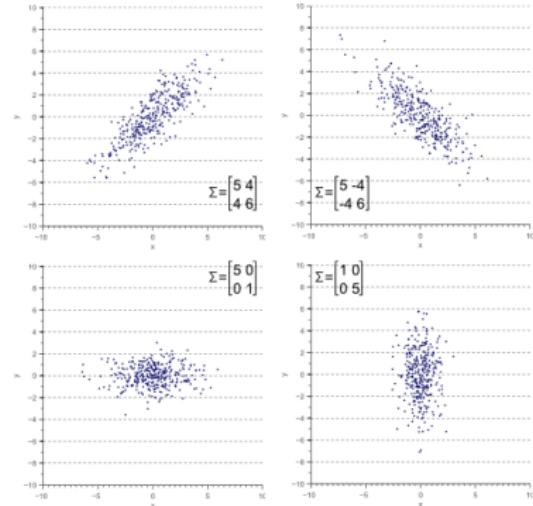
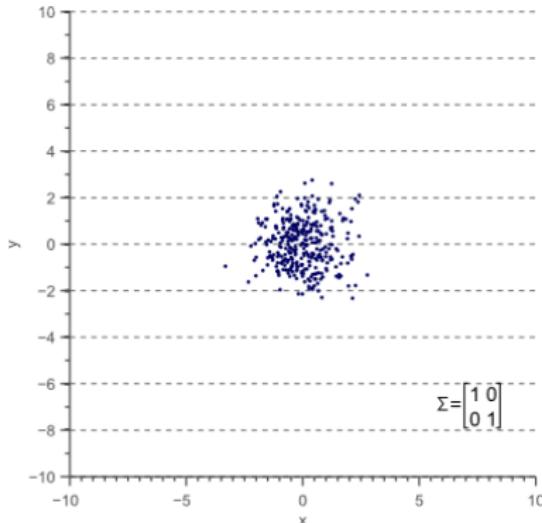
$$Cov(X_j, X_{j'}) = \hat{S}_{jj'} = \frac{1}{n-1} \sum_{i=1}^n X_{ij} X_{ij'}$$

D_S is the $p \times p$ matrix with same diagonal as S

Covariance versus Correlation



Sample Covariance (Correlation) Linear Transformation



Eigenvectors and Eigenvalues

The *eigenvectors* and *eigenvalues* of a linear transformation $\Sigma_{p \times p}$ are orthonormal vectors v and constants λ satisfying

$$\Sigma v = \lambda v$$

$$\begin{bmatrix} 2 & 0 \\ 0 & 1 \end{bmatrix} \begin{bmatrix} 1 \\ 0 \end{bmatrix} = 2 \begin{bmatrix} 1 \\ 0 \end{bmatrix}$$

Eigenvectors and Eigenvalues

The *eigenvectors* and *eigenvalues* of a linear transformation $\Sigma_{p \times p}$ are orthonormal vectors v and constants λ satisfying

$$\Sigma v = \lambda v \quad \text{as opposed to} \quad \Sigma w = w'$$

$$\begin{bmatrix} 2 & 0 \\ 0 & 1 \end{bmatrix} \begin{bmatrix} 1 \\ 0 \end{bmatrix} = 2 \begin{bmatrix} 1 \\ 0 \end{bmatrix} \quad \begin{bmatrix} 2 & 0 \\ 0 & 1 \end{bmatrix} \begin{bmatrix} \frac{1}{\sqrt{2}} \\ \frac{1}{\sqrt{2}} \end{bmatrix} = \begin{bmatrix} \frac{2}{\sqrt{2}} \\ \frac{1}{\sqrt{2}} \end{bmatrix}$$

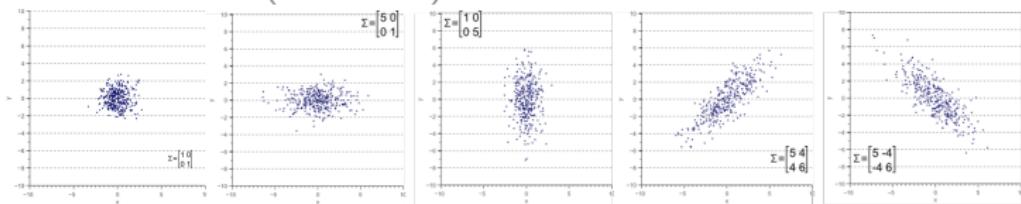
Eigenvectors and Eigenvalues

The *eigenvectors* and *eigenvalues* of a linear transformation $\Sigma_{p \times p}$ are orthonormal vectors v and constants λ satisfying

$$\Sigma v = \lambda v \quad \text{as opposed to} \quad \Sigma w = w'$$

$$\begin{bmatrix} 2 & 0 \\ 0 & 1 \end{bmatrix} \begin{bmatrix} 1 \\ 0 \end{bmatrix} = 2 \begin{bmatrix} 1 \\ 0 \end{bmatrix} \quad \begin{bmatrix} 2 & 0 \\ 0 & 1 \end{bmatrix} \begin{bmatrix} \frac{1}{\sqrt{2}} \\ \frac{1}{\sqrt{2}} \end{bmatrix} = \begin{bmatrix} \frac{2}{\sqrt{2}} \\ \frac{1}{\sqrt{2}} \end{bmatrix}$$

The covariance (correlation) matrix correlates uncorrelated variables



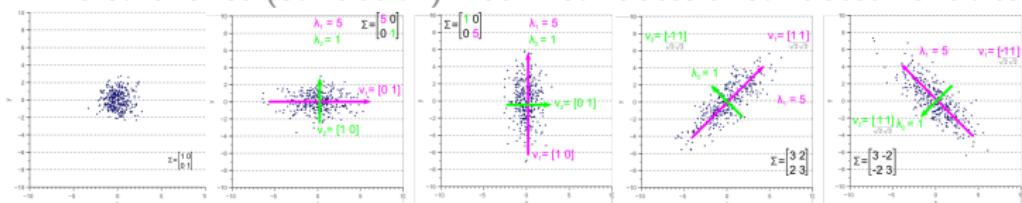
Eigenvectors and Eigenvalues

The *eigenvectors* and *eigenvalues* of a linear transformation $\Sigma_{p \times p}$ are orthonormal vectors v and constants λ satisfying

$$\Sigma v = \lambda v \quad \text{as opposed to} \quad \Sigma w = w'$$

$$\begin{bmatrix} 2 & 0 \\ 0 & 1 \end{bmatrix} \begin{bmatrix} 1 \\ 0 \end{bmatrix} = 2 \begin{bmatrix} 1 \\ 0 \end{bmatrix} \quad \begin{bmatrix} 2 & 0 \\ 0 & 1 \end{bmatrix} \begin{bmatrix} \frac{1}{\sqrt{2}} \\ \frac{1}{\sqrt{2}} \end{bmatrix} = \begin{bmatrix} \frac{2}{\sqrt{2}} \\ \frac{1}{\sqrt{2}} \end{bmatrix}$$

The covariance (correlation) matrix correlates uncorrelated variables



The *eigenvectors* of the covariance (correlation) matrix turn out to be unit vectors v_j sequentially maximizing the variance of scores $z_{ij} = v_j^T X_i$, that are sequentially uncorrelated with previous scores.

i.e. for each j , v_j solve $\underset{v_j}{\operatorname{argmax}} \sum_{j'=1}^j \operatorname{Var}(z_{ij'}) = \sum_{j'=1}^j \lambda_{j'}$ such

that $\operatorname{Cov}(z_{ij}, z_{ij'}) = 0$ for $j' < j$, with λ_j the eigenvalue of v_j .

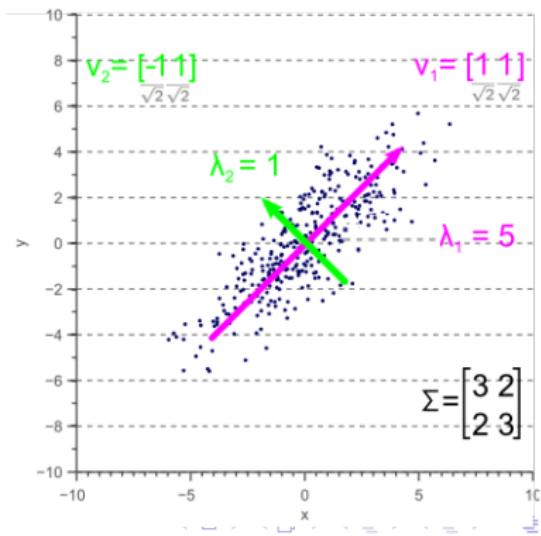
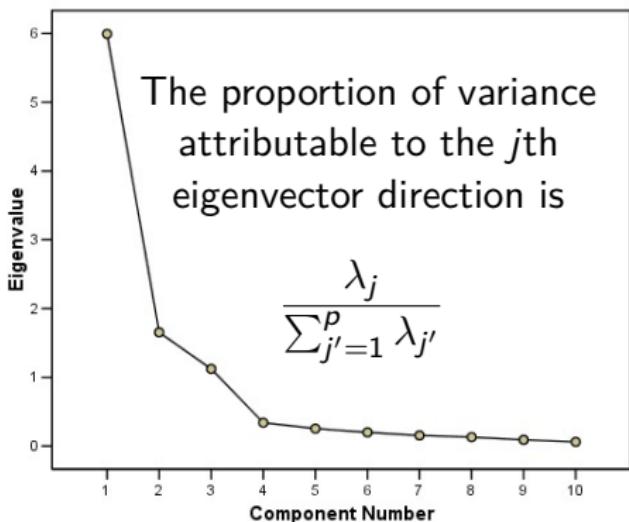
Proportion of Variation Explained

Since our eigenvalues are orthonormal, i.e.,

$$\sum_{k=1}^p v_{kj}^2 = 1 \text{ and } v_j^T v_{j'} = 0 \text{ for } j \neq j' \text{ we have that}$$

$$\sum_{j=1}^p \text{Var}_i(X_{ij}) = \sum_{j=1}^p \text{Var}_i(v_j^T X_{i\cdot}) = \sum_{j=1}^p \text{Var}_i(z_{ij}) = \sum_{j=1}^p \lambda_j$$

i.e., data variance can be partitioned between eigenvector directions



Biasing the “Proportion of Variance Explained” Direction

The eigenvalues point towards the most correlated axes, e.g.

$$\begin{bmatrix} 1 & 0.9 & 0.1 & 0.1 \\ 0.9 & 1 & 0.1 & 0.1 \\ 0.1 & 0.1 & 1 & 0.9 \\ 0.1 & 0.1 & 0.9 & 1 \end{bmatrix} \quad v_1 = 1/\sqrt{\frac{1}{4}} \begin{bmatrix} 1 \\ 1 \\ 1 \\ 1 \end{bmatrix} \quad v_2 = 1/\sqrt{\frac{1}{4}} \begin{bmatrix} -1 \\ -1 \\ 1 \\ 1 \end{bmatrix}$$

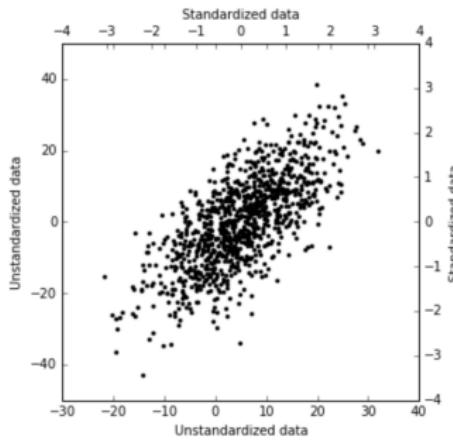
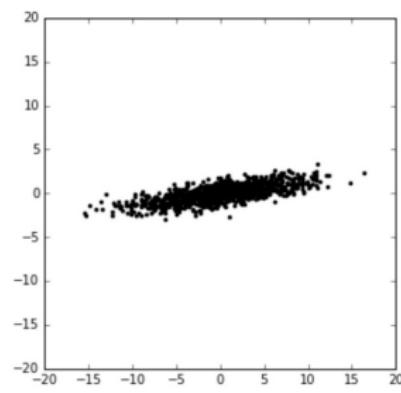
$$\begin{bmatrix} 1 & .6 & .5 \\ .6 & 1 & .4 \\ .5 & .4 & 1 \end{bmatrix} \quad v_1 = \begin{bmatrix} 0.61 \\ 0.58 \\ 0.54 \end{bmatrix}$$

Biasing the “Proportion of Variance Explained” Direction

The eigenvalues point towards the most correlated axes, e.g.

$$\begin{bmatrix} 1 & 0.9 & 0.1 & 0.1 \\ 0.9 & 1 & 0.1 & 0.1 \\ 0.1 & 0.1 & 1 & 0.9 \\ 0.1 & 0.1 & 0.9 & 1 \end{bmatrix} \quad v_1 = 1/\sqrt{\frac{1}{4}} \begin{bmatrix} 1 \\ 1 \\ 1 \\ 1 \end{bmatrix} \quad v_2 = 1/\sqrt{\frac{1}{4}} \begin{bmatrix} -1 \\ -1 \\ 1 \\ 1 \end{bmatrix}$$
$$\begin{bmatrix} 1 & .6 & .5 \\ .6 & 1 & .4 \\ .5 & .4 & 1 \end{bmatrix} \quad v_1 = \begin{bmatrix} 0.61 \\ 0.58 \\ 0.54 \end{bmatrix}$$

The eigenvalues further point towards high variance axes, e.g.

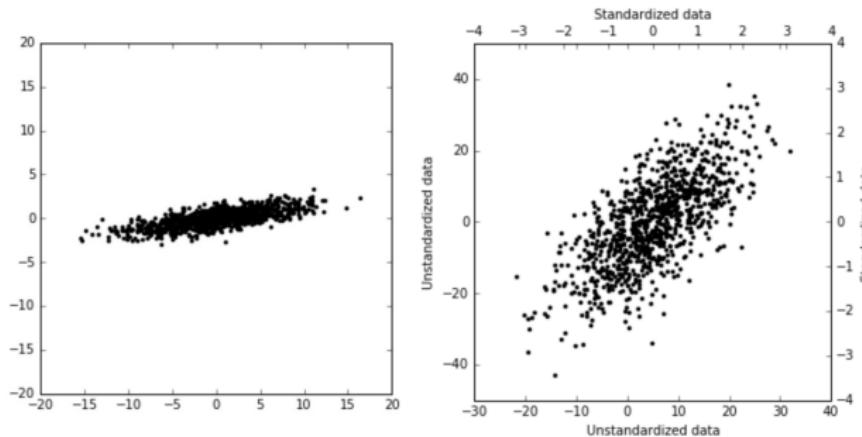


Biasing the “Proportion of Variance Explained” Direction

The eigenvalues point towards the most correlated axes, e.g.

$$\begin{bmatrix} 1 & 0.9 & 0.1 & 0.1 \\ 0.9 & 1 & 0.1 & 0.1 \\ 0.1 & 0.1 & 1 & 0.9 \\ 0.1 & 0.1 & 0.9 & 1 \end{bmatrix} \quad v_1 = 1/\sqrt{\frac{1}{4}} \begin{bmatrix} 1 \\ 1 \\ 1 \\ 1 \end{bmatrix} \quad v_2 = 1/\sqrt{\frac{1}{4}} \begin{bmatrix} -1 \\ -1 \\ 1 \\ 1 \end{bmatrix}$$
$$\begin{bmatrix} 1 & .6 & .5 \\ .6 & 1 & .4 \\ .5 & .4 & 1 \end{bmatrix} \quad v_1 = \begin{bmatrix} 0.61 \\ 0.58 \\ 0.54 \end{bmatrix}$$

The eigenvalues further point towards high variance axes, e.g.



When is $\sum_{j=1}^p \lambda_j = p$?

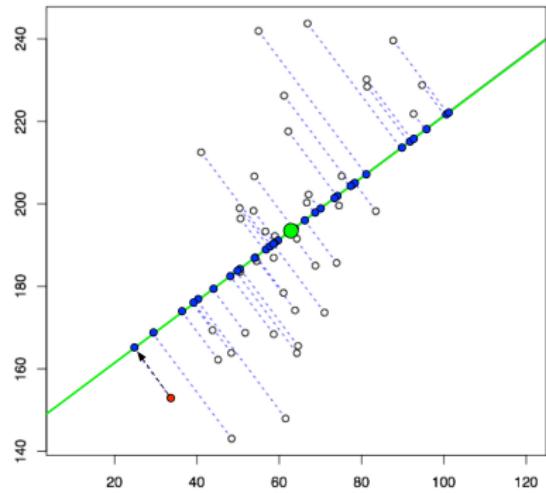
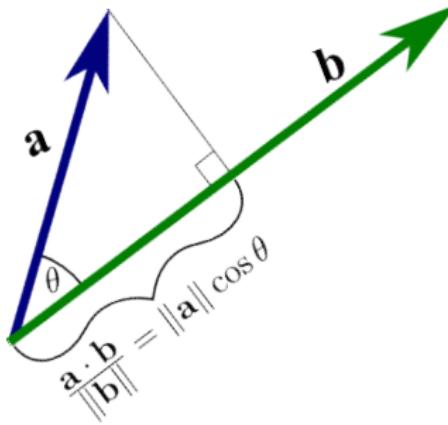
Mapping onto the directions of “Greatest Variation”

The orthogonally projection of X_i onto v_j is

$$z_{ij} = \frac{v_j}{\|v_j\|} \cdot X_i$$

For unit *eigenvectors* v_j , this is simply the dot product

$$z_{ij} = v_j \cdot X_i$$



Principal Components Analysis (PCA)

- ▶ PCA projects \mathbf{X} onto the \mathbf{X} 's covariance matrix eigenvectors
- ▶ An eigenvector projection is called a principal component (PC)
- ▶ The 1^{st} PC is the projection onto the eigenvector with the *largest eigenvalue*, the 2^{nd} PC is the projection onto the eigenvector with the *second largest eigenvalue*, etc.

Principal Components Analysis (PCA)

- ▶ PCA projects \mathbf{X} onto the \mathbf{X} 's covariance matrix eigenvectors
- ▶ An eigenvector projection is called a principal component (PC)
- ▶ The 1st PC is the projection onto the eigenvector with the *largest eigenvalue*, the 2nd PC is the projection onto the eigenvector with the *second largest eigenvalue*, etc.
- ▶ The *Principal Components \mathbf{z}* are an *orthogonal rotation* of \mathbf{X}

Feature	X_1	X_2	\cdots	X_p
n=1				
2				
:				
n				

$\Sigma \rightarrow$

PC	z_1	z_2	\cdots	z_p
n=1				
2				
:				
n				

Principal Components Analysis (PCA)

- ▶ PCA projects \mathbf{X} onto the \mathbf{X} 's covariance matrix eigenvectors
- ▶ An eigenvector projection is called a principal component (PC)
- ▶ The 1st PC is the projection onto the eigenvector with the *largest eigenvalue*, the 2nd PC is the projection onto the eigenvector with the *second largest eigenvalue*, etc.
- ▶ The *Principal Components \mathbf{z}* are an *orthogonal rotation* of \mathbf{X}

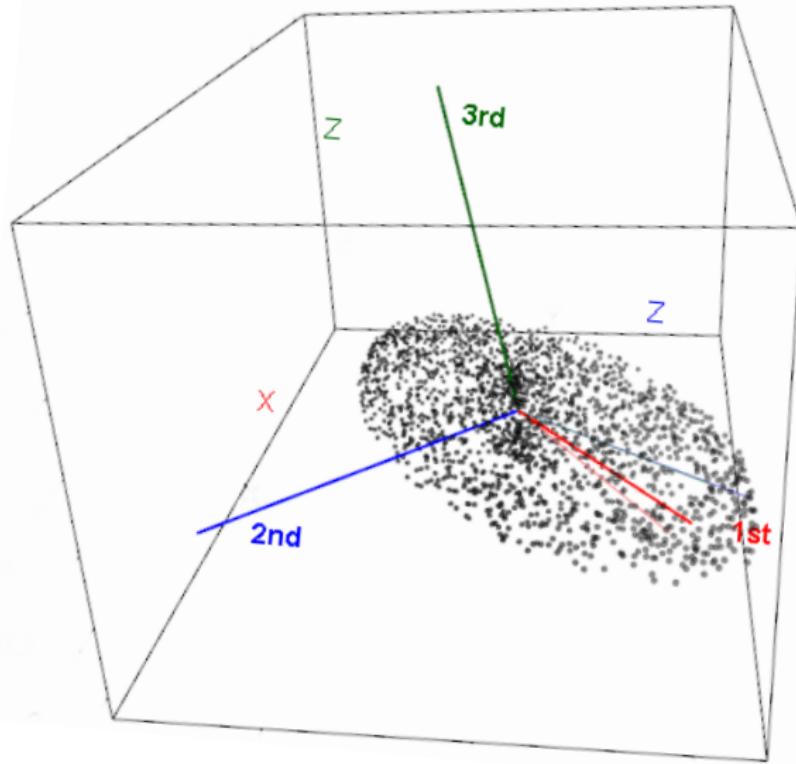
Feature	X_1	X_2	\cdots	X_p
n=1				
2				
:				
n				

Σ

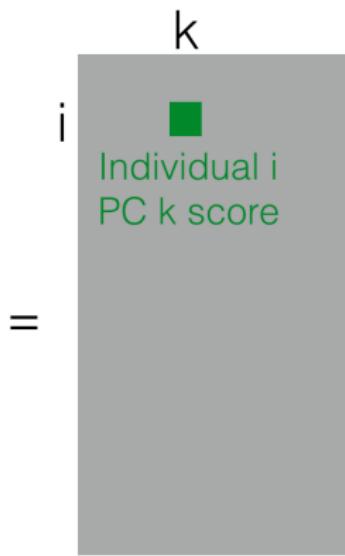
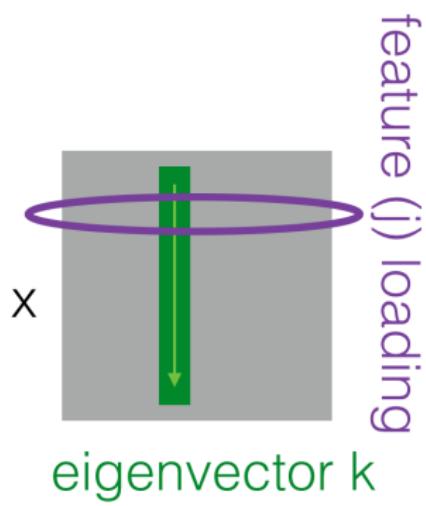
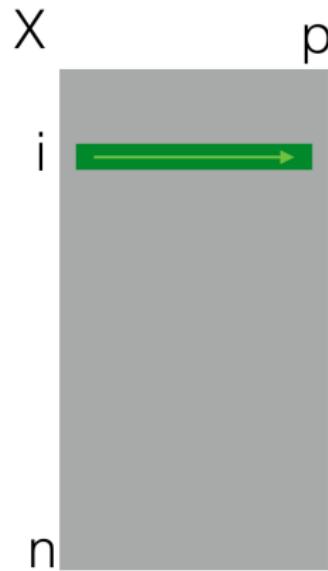
PC	z_1	z_2	\cdots	z_p
n=1				
2				
:				
n				

- ▶ PCA routinely centers and scales $X_{\cdot j}$ in order to capture correlation between features rather than feature scales

PCA: Transforms Features To Uncorrelated Features



Mechanics and Terms



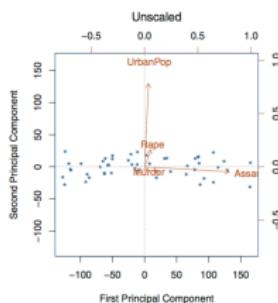
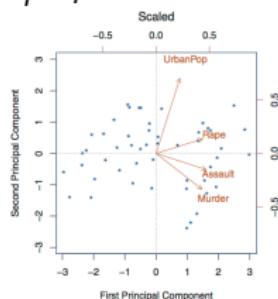
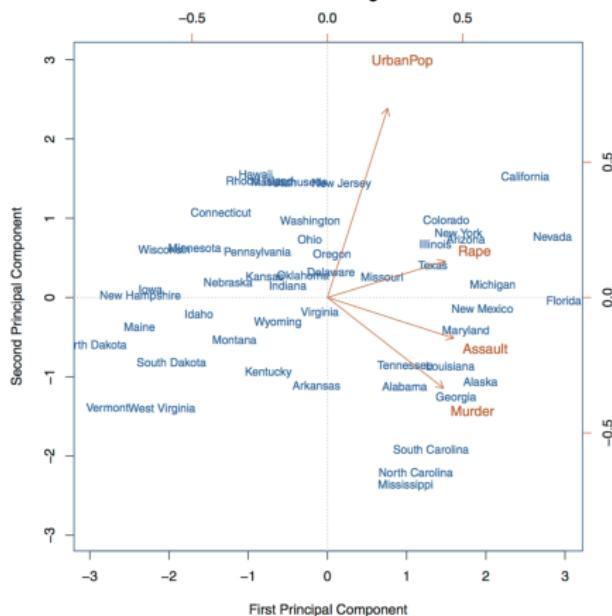
Biplots: look at the axes

- ▶ **Vector:** k^{th} versus k'^{th} eigenvector's j^{th} feature loading:

v_{kj} versus $v_{k'j}$

- ▶ **Point:** observation i 's j^{th} versus j'^{th} PC scores:

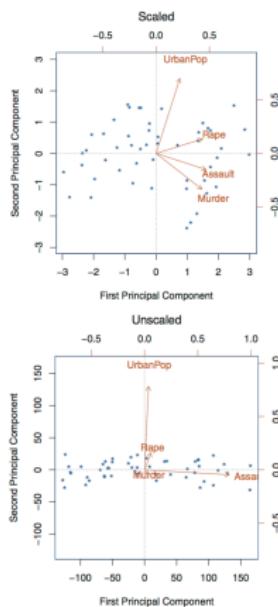
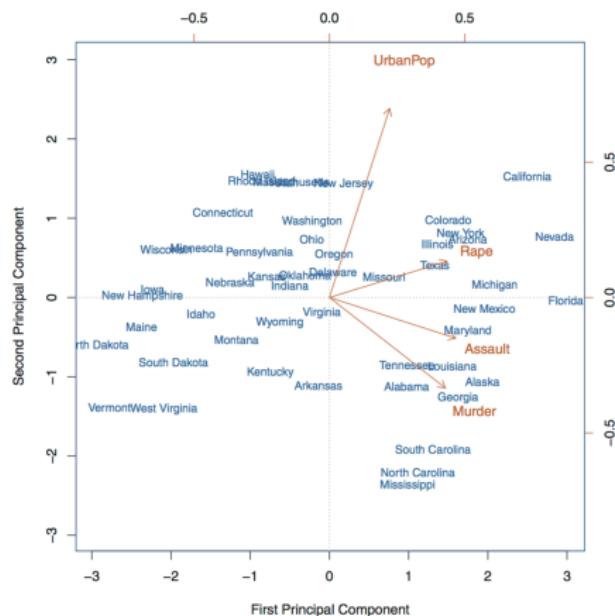
$$z_{ji} = v_j' X_i \text{ versus } z_{j'i} = v_{i'}' X_i$$



Biplots: look at the axes

- ▶ **Vector:** which features the PCs point towards

- ▶ **Point:** which data points are driven by which PCs



Dimensionality Reduction?

2016 Volvo S60 T5 Inscription  Great Deal 52 days on CarGurus Price: \$30,250 \$563/mo est. Mileage: 6,137 mi Location: Austin, TX 3 mi Dealer rating: ★★★★☆ Save	2015 Volvo S60 2015.5 T5 Premier  Good Deal 9 days on CarGurus Price: \$23,995 \$447/mo est. Mileage: 17,141 mi Location: Austin, TX 5 mi Dealer rating: ★★★★☆ Save
2014 Volvo S60 T5  Great Deal 83 days on CarGurus Price: \$16,995 \$317/mo est. Mileage: 48,058 mi Location: Austin, TX 5 mi Dealer rating: ★★★★☆ Save	2015 Volvo S60 T5 AWD  Good Deal 43 days on CarGurus Price: \$18,250 \$340/mo est. Mileage: 50,936 mi Location: Austin, TX 3 mi Dealer rating: ★★★★☆ Save
2013 Volvo S60 T5  Great Deal 181 days on CarGurus Price: \$14,970 \$279/mo est. Mileage: 36,626 mi Location: Austin, TX 3 mi Dealer rating: ★★★★☆ Save	2015 Volvo S60 T5 Premier Plus  Fair Deal 10 days on CarGurus Price: \$23,495 \$437/mo est. Mileage: 26,927 mi Location: Austin, TX 5 mi Dealer rating: ★★★★☆ Save
2012 Volvo S60 T5  Good Deal 7 days on CarGurus Price: \$10,000 \$186/mo est. Mileage: 98,917 mi Location: San Marcos, TX 35 mi Dealer rating: ★★★★☆ Save	2016 Volvo S60 T5  Great Deal 61 days on CarGurus Price: \$26,000 \$484/mo est. Mileage: 7,405 mi Location: Austin, TX 3 mi Dealer rating: ★★★★☆ Save
2015 Volvo S60 2015.5 T5 Premier  Good Deal 9 days on CarGurus Price: \$23,995 \$447/mo est. Mileage: 17,141 mi Location: Austin, TX 5 mi Dealer rating: ★★★★☆ Save	2014 Volvo S60 T5  Great Deal 25 days on CarGurus Price: \$20,250 \$377/mo est. Mileage: 11,140 mi Location: Austin, TX 3 mi Dealer rating: ★★★★☆ Save

Dimensionality Reduction?

 CARVANA	2016 Volvo S60 T5 Inscription  Great Deal 52 days on CarGurus Price: \$30,250 \$563/mo est. Mileage: 6,137 mi Location: Austin, TX 3 mi Dealer rating: ★★★★☆		 CARVANA	2015 Volvo S60 2015.5 T5 Premier  Good Deal 9 days on CarGurus Price: \$23,995 \$447/mo est. Mileage: 17,141 mi Location: Austin, TX 5 mi Dealer rating: ★★★★☆	
 CARVANA	2014 Volvo S60 T5  Great Deal 83 days on CarGurus Price: \$16,995 \$317/mo est. Mileage: 48,058 mi Location: Austin, TX 5 mi Dealer rating: ★★★★☆		 CARVANA	2015 Volvo S60 T5 AWD  Good Deal 43 days on CarGurus Price: \$18,250 \$340/mo est. Mileage: 50,936 mi Location: Austin, TX 3 mi Dealer rating: ★★★★☆	
 CARVANA	2013 Volvo S60 T5  Great Deal 181 days on CarGurus Price: \$14,970 \$279/mo est. Mileage: 36,626 mi Location: Austin, TX 3 mi Dealer rating: ★★★★☆		 CARVANA	2015 Volvo S60 T5 Premier Plus  Fair Deal 10 days on CarGurus Price: \$23,495 \$437/mo est. Mileage: 26,927 mi Location: Austin, TX 5 mi Dealer rating: ★★★★☆	
 CARVANA	2012 Volvo S60 T5  Good Deal 7 days on CarGurus Price: \$10,000 \$186/mo est. Mileage: 96,917 mi Location: San Marcos, TX 35 mi Dealer rating: ★★★★☆		 CARVANA	2016 Volvo S60 T5  Great Deal 61 days on CarGurus Price: \$26,000 \$484/mo est. Mileage: 7,405 mi Location: Austin, TX 3 mi Dealer rating: ★★★★☆	
 CARVANA	2015 Volvo S60 2015.5 T5 Premier  Good Deal 9 days on CarGurus Price: \$23,995 \$447/mo est. Mileage: 17,141 mi Location: Austin, TX 5 mi Dealer rating: ★★★★☆		 CARVANA	2014 Volvo S60 T5  Great Deal 25 days on CarGurus Price: \$20,250 \$377/mo est. Mileage: 11,140 mi Location: Austin, TX 3 mi Dealer rating: ★★★★☆	

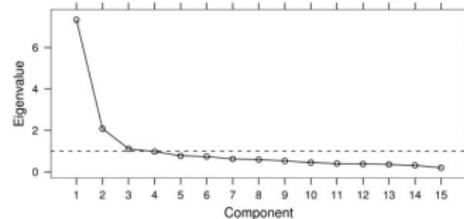
Dimensionality Reduction?

How does e.g. KNN feel about this many dimensions? Curse of...

 2016 Volvo S60 T5 Inscription Great Deal 52 days on CarGurus Price: \$30,250 \$563/mo est. Mileage: 6,137 mi Location: Austin, TX 3 mi Dealer rating: ★★★★☆	 2015 Volvo S60 2015.5 T5 Premier Good Deal 9 days on CarGurus Price: \$23,995 \$447/mo est. Mileage: 17,141 mi Location: Austin, TX 5 mi Dealer rating: ★★★★☆
 2014 Volvo S60 T5 Great Deal 83 days on CarGurus Price: \$16,995 \$317/mo est. Mileage: 48,058 mi Location: Austin, TX 5 mi Dealer rating: ★★★★☆	 2015 Volvo S60 T5 AWD Good Deal 43 days on CarGurus Price: \$18,250 \$340/mo est. Mileage: 50,936 mi Location: Austin, TX 3 mi Dealer rating: ★★★★☆
 2013 Volvo S60 T5 Great Deal 181 days on CarGurus Price: \$14,970 \$279/mo est. Mileage: 36,626 mi Location: Austin, TX 3 mi Dealer rating: ★★★★☆	 2015 Volvo S60 T5 Premier Plus Fair Deal 10 days on CarGurus Price: \$23,495 \$437/mo est. Mileage: 26,927 mi Location: Austin, TX 5 mi Dealer rating: ★★★★☆
 2012 Volvo S60 T5 Good Deal 7 days on CarGurus Price: \$10,000 \$186/mo est. Mileage: 96,917 mi Location: San Marcos, TX 35 mi Dealer rating: ★★★★☆	 2016 Volvo S60 T5 Great Deal 61 days on CarGurus Price: \$26,000 \$484/mo est. Mileage: 7,405 mi Location: Austin, TX 3 mi Dealer rating: ★★★★☆
 2015 Volvo S60 2015.5 T5 Premier Good Deal 9 days on CarGurus Price: \$23,995 \$447/mo est. Mileage: 17,141 mi Location: Austin, TX 5 mi Dealer rating: ★★★★☆	 2014 Volvo S60 T5 Great Deal 25 days on CarGurus Price: \$20,250 \$377/mo est. Mileage: 11,140 mi Location: Austin, TX 3 mi Dealer rating: ★★★★☆

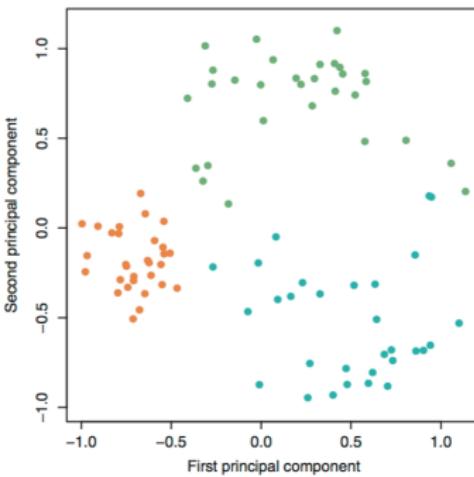
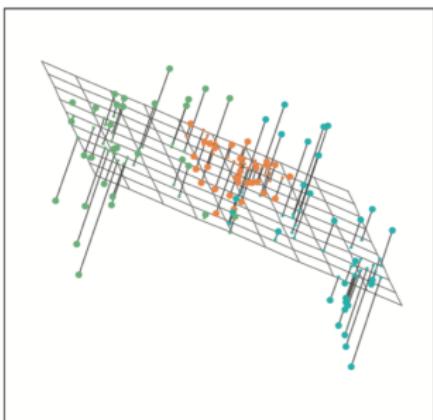
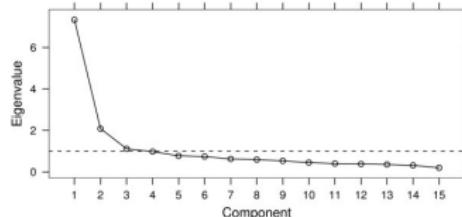
Dimensionality Reduction

$$X_i \underset{p \times 1}{\approx} \sum_{j=1}^{q < p} v_j \underset{p \times 1}{z_{ij}}$$



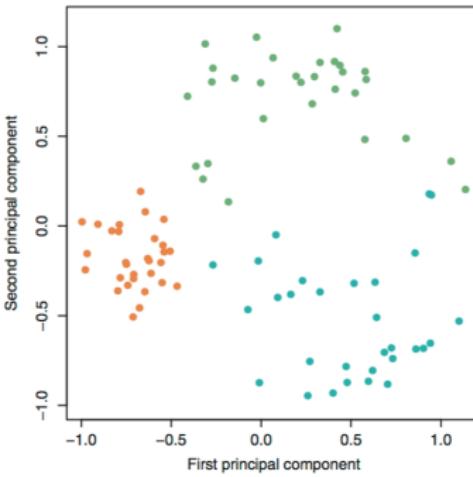
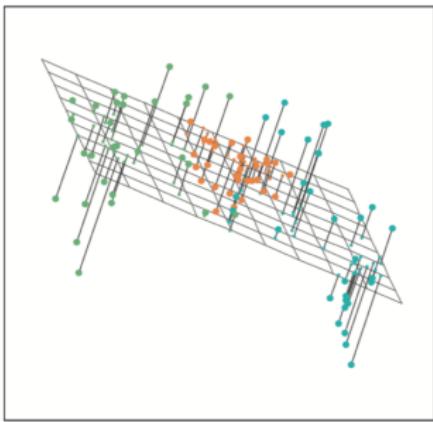
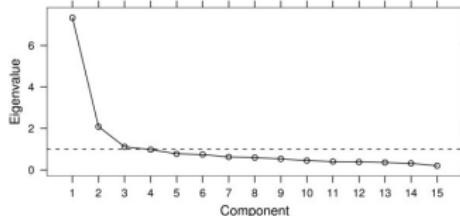
Dimensionality Reduction

$$X_i \underset{p \times 1}{\approx} \sum_{j=1}^{q < p} v_j \underset{p \times 1}{z_{ij}}$$



Dimensionality Reduction

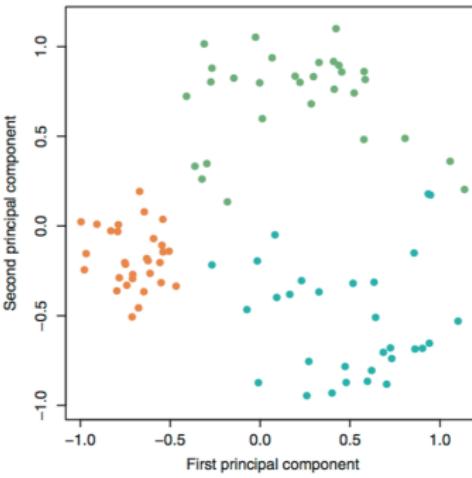
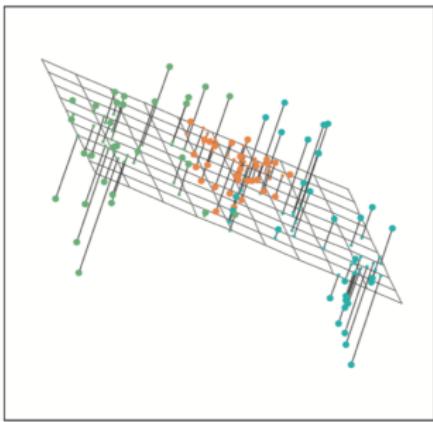
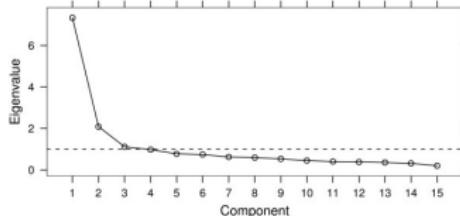
$$X_i \underset{p \times 1}{\approx} \sum_{j=1}^{q < p} v_j z_{ij} \underset{p \times 1}{}$$



Compresses data into uncorrelated variation & maybe drops noise
(Feature Engineering?) (Curse of Dimensionality?)

Dimensionality Reduction

$$\underset{p \times 1}{X_i} \approx \sum_{j=1}^{q < p} \underset{p \times 1}{v_j} z_{ij}$$



Compresses data into uncorrelated variation & maybe drops noise
(Feature Engineering?) (Curse of Dimensionality?)

The number of nonzero λ_j 's equals the rank of the covariance matrix $\mathbf{X}^T \mathbf{X}$ which itself equals the rank of the data matrix \mathbf{X}

Dimensionality Reduction

Feature	X_1	X_2	\dots	X_p
$n=1$				
2				
\vdots				
n				

$\Sigma \rightarrow$

PC	z_1	z_2	\dots	z_p
$n=1$				
2				
\vdots				
n				

Dimensionality Reduction

Feature	X_1	X_2	\cdots	X_p
$n=1$				
2				
\vdots				
n				

Σ

PC	z_1	z_2	\cdots	z_p
$n=1$				
2				
\vdots				
n				

Feature	X_1	X_2	\cdots	X_p
$n=1$				
2				
\vdots				
n				

Σ

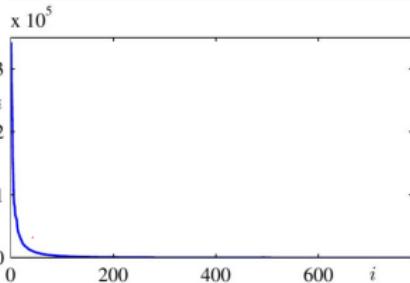
PC	z_1	\cdots
$n=1$		
2		
\vdots		
n		

Dimensionality Reduction Examples

00000000000000000000
111111111111111111
222222222222222222
333333333333333333
444444444444444444
555555555555555555
666666666666666666
777777777777777777
888888888888888888
999999999999999999

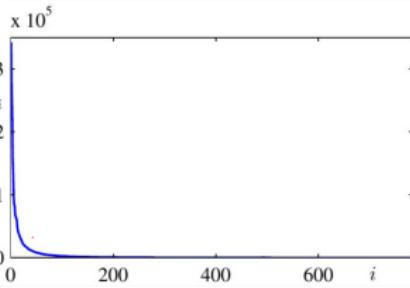
Dimensionality Reduction Examples

```
0 0 0 0 0 0 0 0 0 0 0 0 0 0 0 0 0 0 0 0 0 0 0 0  
/ / / / / / / / / / / / / / / / / / / / / / / /  
2 2 2 2 2 2 2 2 2 2 2 2 2 2 2 2 2 2 2 2 2 2  
8 3 3 3 3 3 3 3 3 3 3 3 3 3 3 3 3 3 3 3 3 3  
4 4 4 4 4 4 4 4 4 4 4 4 4 4 4 4 4 4 4 4 4 4  
5 5 5 5 5 5 5 5 5 5 5 5 5 5 5 5 5 5 5 5 5 5  
6 6 6 6 6 6 6 6 6 6 6 6 6 6 6 6 6 6 6 6 6 6  
7 7 7 7 7 7 7 7 7 7 7 7 7 7 7 7 7 7 7 7 7  
8 8 8 8 8 8 8 8 8 8 8 8 8 8 8 8 8 8 8 8 8 4  
9 9 9 9 9 9 9 9 9 9 9 9 9 9 9 9 9 9 9 9 9 4
```



Dimensionality Reduction Examples

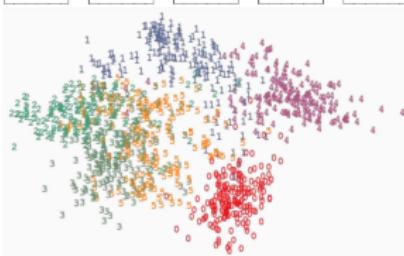
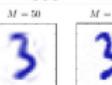
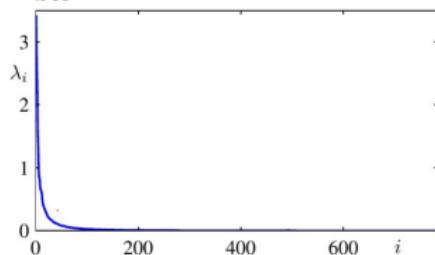
```
0 0 0 0 0 0 0 0 0 0 0 0 0 0 0 0 0 0 0 0 0 0 0 0  
/ / / / / / / / / / / / / / / / / / / / / / / /  
2 2 2 2 2 2 2 2 2 2 2 2 2 2 2 2 2 2 2 2 2 2 2  
8 3 3 3 3 3 3 3 3 3 3 3 3 3 3 3 3 3 3 3 3 3  
4 4 4 4 4 4 4 4 4 4 4 4 4 4 4 4 4 4 4 4 4 4  
5 5 5 5 5 5 5 5 5 5 5 5 5 5 5 5 5 5 5 5 5 5  
6 6 6 6 6 6 6 6 6 6 6 6 6 6 6 6 6 6 6 6 6 6  
7 7 7 7 7 7 7 7 7 7 7 7 7 7 7 7 7 7 7 7 7 7  
8 8 8 8 8 8 8 8 8 8 8 8 8 8 8 8 8 8 8 8 8 8  
9 9 9 9 9 9 9 9 9 9 9 9 9 9 9 9 9 9 9 9 9 9
```



Dimensionality Reduction Examples

0 0 0 0 0 0 0 0 0 0 0 0 0 0 0 0 0 0
1 1 1 1 1 1 1 1 1 1 1 1 1 1 1 1 1 1
2 2 2 2 2 2 2 2 2 2 2 2 2 2 2 2 2 2
3 3 3 3 3 3 3 3 3 3 3 3 3 3 3 3 3 3
4 4 4 4 4 4 4 4 4 4 4 4 4 4 4 4 4 4
5 5 5 5 5 5 5 5 5 5 5 5 5 5 5 5 5 5
6 6 6 6 6 6 6 6 6 6 6 6 6 6 6 6 6 6
7 7 7 7 7 7 7 7 7 7 7 7 7 7 7 7 7 7
8 8 8 8 8 8 8 8 8 8 8 8 8 8 8 8 8 8
9 9 9 9 9 9 9 9 9 9 9 9 9 9 9 9 9 9

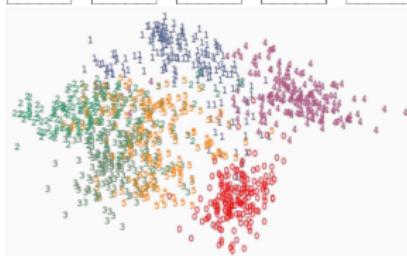
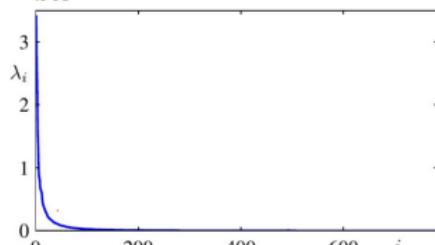
$\times 10^5$



Dimensionality Reduction Examples

0
1
2
3
4
5
6
7
8 4
9 4

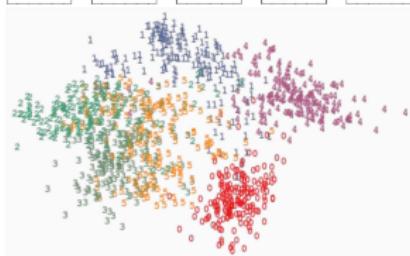
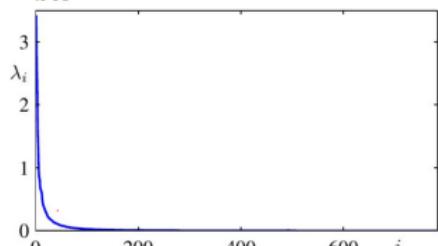
$\times 10^5$



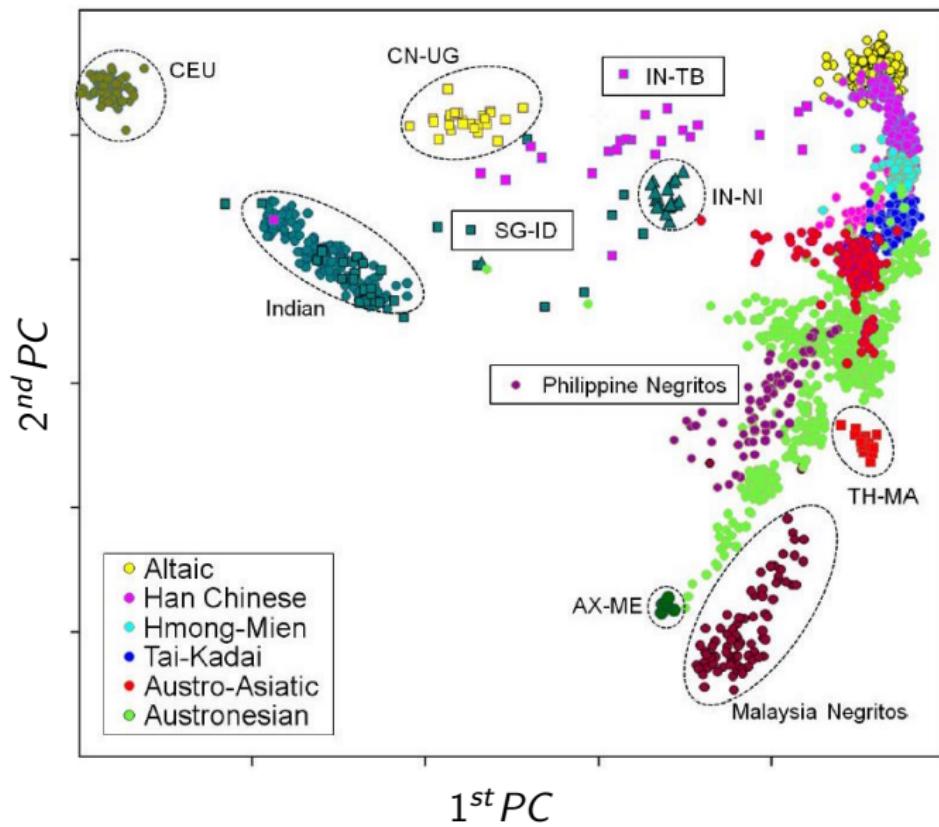
Dimensionality Reduction Examples

0
1
2
3
4
5
6
7
8 4
9 4

$\times 10^5$



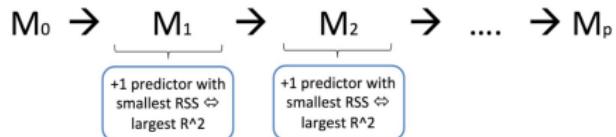
Visualization (for EDA purposes)



How else could we do dimension reduction?

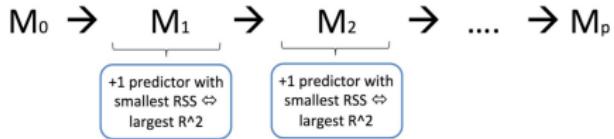
How else could we do dimension reduction?

- ▶ Stepwise Selection

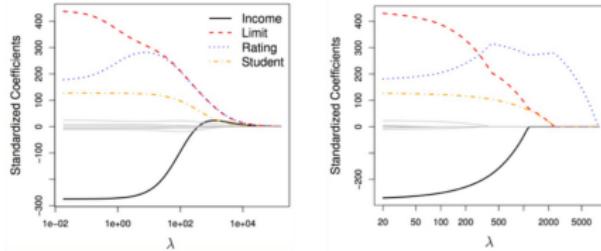


How else could we do dimension reduction?

- ▶ Stepwise Selection

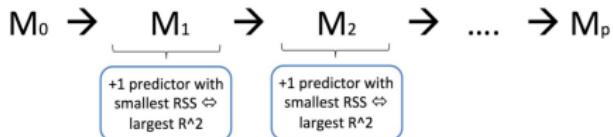


- ▶ Lasso

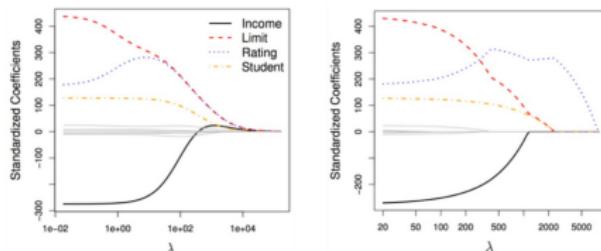


How else could we do dimension reduction?

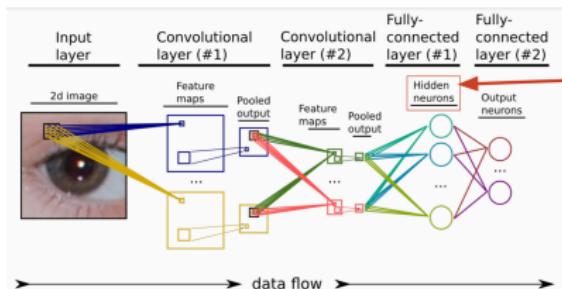
- ▶ Stepwise Selection



- ▶ Lasso

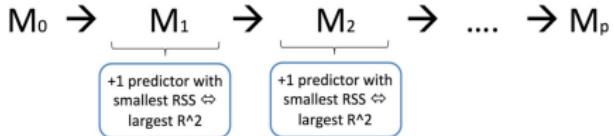


- ▶ Neural Networks

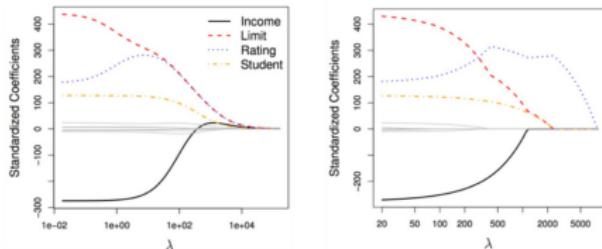


How else could we do dimension reduction?

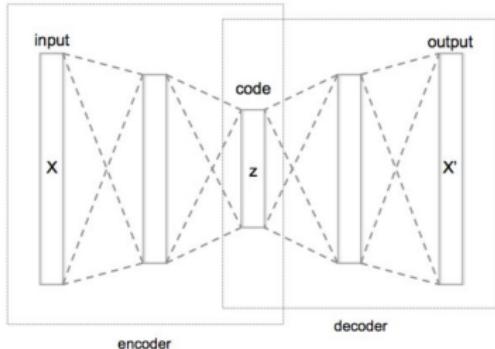
- ▶ Stepwise Selection



- ▶ Lasso

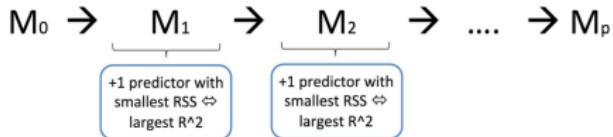


- ▶ Neural Networks

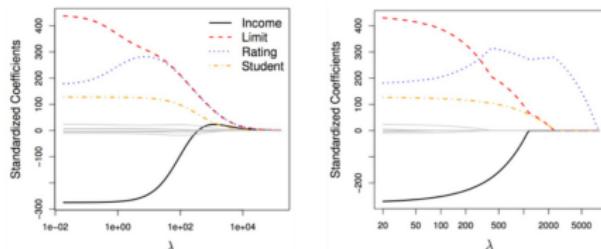


How else could we do dimension reduction?

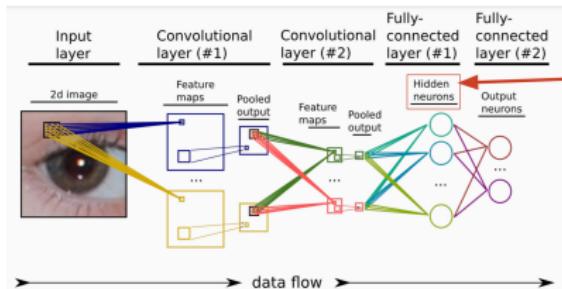
- ▶ Stepwise Selection



- ▶ Lasso

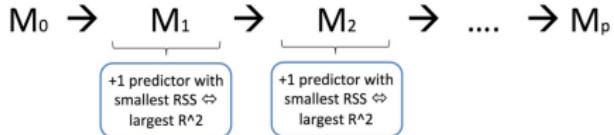


- ▶ Neural Networks

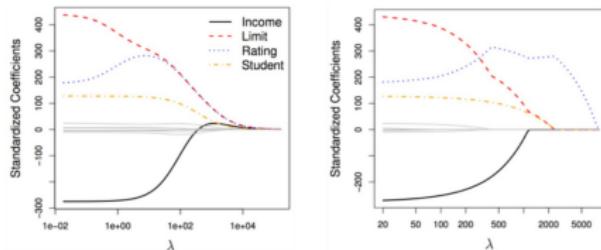


How else could we do dimension reduction?

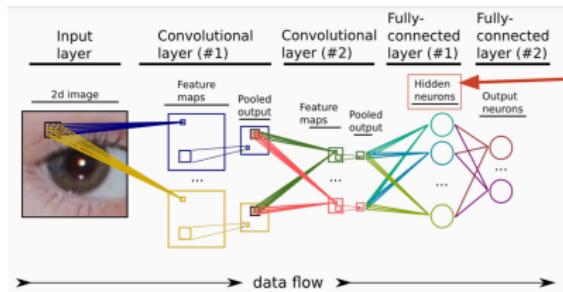
- ▶ Stepwise Selection



- ▶ Lasso



- ▶ Neural Networks



It's fun to stay at the

- ▶ PCA, eh?

Singular Value Decomposition (SVD)

$$\text{Data} \quad n \times p = U \quad n \times (q < n) \quad (q < n) \times (q < p) \quad (q < p) \times p$$
$$\Sigma \quad V^T$$

- ▶ Data: “standarized” \mathbf{X}
- ▶ U : Standardized principal component scores
- ▶ Σ : $\sqrt{\text{eigenvalues}}$ on diagonal, 0's off diagonal
- ▶ V^T : Matrix whose rows are the eigenvectors of $\mathbf{X}^T \mathbf{X}$

$$X_{ij'} \approx \sum_{j=1}^{q < p} U_{ij} \sqrt{\lambda_j} V_{jj'}^T$$

$$X_i \approx \sum_{j=1}^{q < p} v_j z_{ij}$$

Latent Factors

$$\text{Data} = \mathbf{U} \Sigma \mathbf{V}^T$$

- ▶ Data is the “phenotypes” (features)
- ▶ U is the “genes” (factors) driving “phenotypes” (features)
- ▶ V^T is the “way” “genes” drive “phenotypes”
- ▶ Σ controls the strength each “genes” influence

Latent Factors

$$\text{Data} = \mathbf{U} \Sigma \mathbf{V}^T$$

- ▶ Data is the “phenotypes” (features)
 - ▶ U is the “genes” (factors) driving “phenotypes” (features)
 - ▶ V^T is the “way” “genes” drive “phenotypes”
 - ▶ Σ controls the strength each “genes” influence
-
- ▶ For some “high dimensional” data, e.g., audio, video, these *latent factors* are the meaningful features associating with Y 's

Latent Factors

$$\text{Data} = \mathbf{U} \Sigma \mathbf{v}^T$$

The diagram illustrates the matrix factorization of a data matrix. On the left, a large gray rectangular box is labeled "Data". An equals sign follows it. To the right of the equals sign is another gray rectangular box divided vertically into two sections: a dark gray left section and a light gray right section, labeled "U". To the right of "U" is a smaller gray rectangular box divided horizontally into two sections: a dark gray top section and a light gray bottom section, labeled "Σ". To the right of "Σ" is a final gray rectangular box divided vertically into two sections: a dark gray left section and a light gray right section, labeled "v^T".

- ▶ When $p \gg n$, $\mathbf{X}^T \mathbf{X}$ is huge
- ▶ SVD doesn't require computation of $\mathbf{X}^T \mathbf{X}$ (i.e., \mathbf{S})
- ▶ SVD is a faster, more stable algorithm (than PCA)

Bonus: Factor Analysis – genes and phenotypes

$$X_{ij} - \mu_j = \sum_{k=1}^q \beta_{jk} F_{ik} + \epsilon_{ij}$$

	What	Index	Variable	Note
n	samples	$i = 1, \dots, n$		
p	features/sample	$j = 1, \dots, p$	X_{ij}	$E[X_{ij}] = \mu_j$
p	errors/sample	$j = 1, \dots, p$	ϵ_{ij}	$\stackrel{iid}{\sim} N(0, \sigma_j)$
q	factors/sample	$k = 1, \dots, q$	F_{ik}	$\stackrel{iid}{\sim} N(0, 1)$
$p \times q$	factor loadings		β_{jk}	Constants

So for each sample i we have p equations

$$X_{i1} - \mu_1 = \beta_{11} F_{i1} + \dots + \beta_{1k} F_{ik} + \epsilon_{i1}$$

$$X_{i2} - \mu_2 = \beta_{21} F_{i1} + \dots + \beta_{2k} F_{ik} + \epsilon_{i2}$$

⋮

$$X_{ip} - \mu_p = \beta_{p1} F_{i1} + \dots + \beta_{pk} F_{ik} + \epsilon_{ip}$$

Bonus: Factor Analysis

If there are fewer factors F_{ik} than features X_{ij} , then we have a reduced representation of the information in the features.

For all j , $X_{ij} - \mu_j = \sum_{k=1}^q \beta_{jk} F_{ik} + \epsilon_{ij}$ can be written as

$$\mathbf{X}_i - \boldsymbol{\mu} = \boldsymbol{\beta}\mathbf{F} + \boldsymbol{\epsilon}$$

Given X_{ij} 's, what are β_{jk} , F_{ik} , and ϵ_{ij} ?

$$\begin{aligned} \text{Cov}(\mathbf{X}) &= \text{Cov}(\boldsymbol{\beta}\mathbf{F} + \boldsymbol{\epsilon}) \\ &= \boldsymbol{\beta}\boldsymbol{\beta}^T + \text{Cov}(\boldsymbol{\epsilon}) \end{aligned}$$

and the off diagonal entries of $\text{Cov}(\boldsymbol{\epsilon})$ should be zero, so we look for factor loadings $\boldsymbol{\beta}$ to minimize the off diagonal entries of $\text{Cov}(\boldsymbol{\epsilon})$.

Given $\boldsymbol{\beta}$, \mathbf{F} and $\boldsymbol{\epsilon}$ can be solved for. The $\boldsymbol{\beta}$ can be adjusted to account for rotations in \mathbf{F} that allow for more interpretability.

Bonus: Independent Component Analysis (ICA)

ICA starts out the same way as factor analysis, i.e.,

$$X_{ij} - \mu_j = \sum_{k=1}^q \beta_{jk} F_{ik}$$

but often omits an error term ϵ as it's an effective procedure w/o it

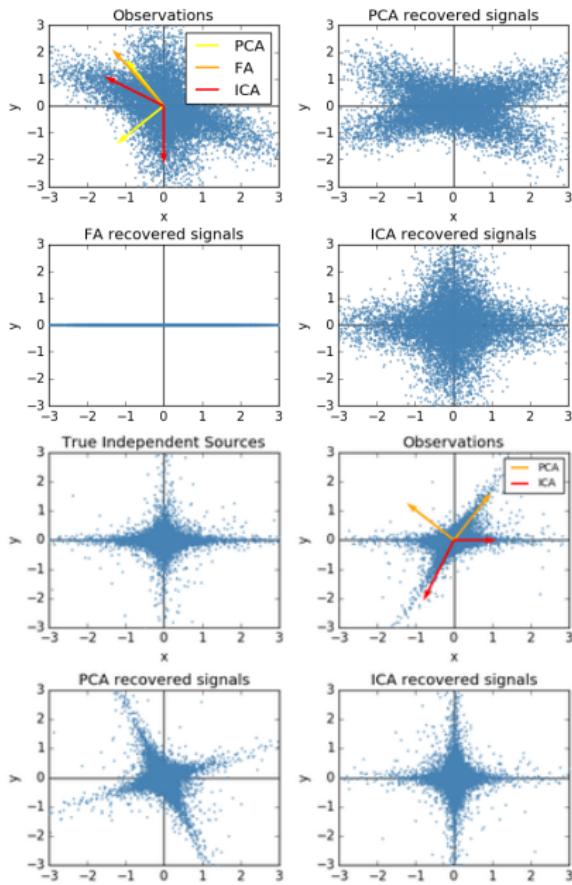
Only independence of F_{ik} 's is assumed – not marginal normality; in fact, the components are chosen sequentially based on *extremeness from normality* (as characterized by kurtosis, their 4th-moment*).

ICA imagines each F_{ik} gives out independent signals picked up to varying degrees by different X_{ij} (i.e., *the cocktail party problem*).

* 1st-moment is mean, 2nd-moment is variance, 3rd-moment is skew.

Kurtosis is a measure of "tail thickness," i.e., the extremeness of points that can be generated by the distribution in question.

Bonus: Graphically



Factor analysis finds an underlying uncorrelated lower dimension multivariate normal distribution that, *plus noise*, generates the observed higher dimension data cloud. E.g., 3D data generated from a 2D uncorrelated multivariate normal distribution:

$$X_{i1} = 1 \cdot F_{i1} + 0 \cdot F_{i2} + \epsilon_{i1}$$

$$X_{i2} = 1 \cdot F_{i1} + 0 \cdot F_{i2} + \epsilon_{i2}$$

$$X_{i3} = 0 \cdot F_{i1} + 1 \cdot F_{i2} + \epsilon_{i3}$$

ICA finds a basis representation of the data cloud such that the coordinates of the representation are independent of each other.

Bonus: PCA vs. Factor Analysis vs. ICA

The eigenvectors of the feature covariance matrix are an orthogonal basis whose axes capture the directions of maximum variation in the data. The eigenvalues are the relative amounts of variation captured along these axes. Transformation of the data to this coordinate system (i.e., *principal components*) is called PCA.

PCA is not a generative model like FA – PCA is purely a representative tool used to describe the data; nonetheless, both methods find orthogonal basis data representations; moreover, the (reduced representation) bases reconstruct the covariance matrix (PCA) or its off-diagonal entries (FA) so for a sufficient number, q , of factors and principal components and a relatively uninfluential diagonal, PCA and FA give similar representations. From here FA latent factors can be “rotated” for interpretation. E.g., a latent factor could be lined up with an ICA component rather than a PCA variance direction. Unlike ICA, however, the collective set of factor analysis factors remain orthogonal.

In ICA the underlying basis *is not* orthogonal. This is because rather than targeting an *uncorrelated* representation like FA and PCA, ICA targets an independent representation. Interestingly, the detection of such a representation requires *exactly* that the underlying representations *not be normal* – the opposite of the FA assumption. Independent signals need not be orthogonal features... and interestingly these are found in ICA by simply maximizing the fourth (kurtosis) rather than the second (variance) moment.

Some Scratch Work

$$\text{Var}(X_1) = 3$$

$$\text{Var}(X_2) = 3$$

$$\text{Cov}(X_1, X_2) = 2$$

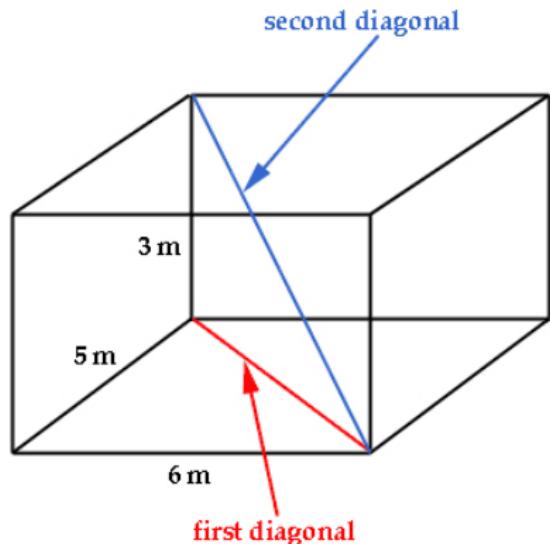
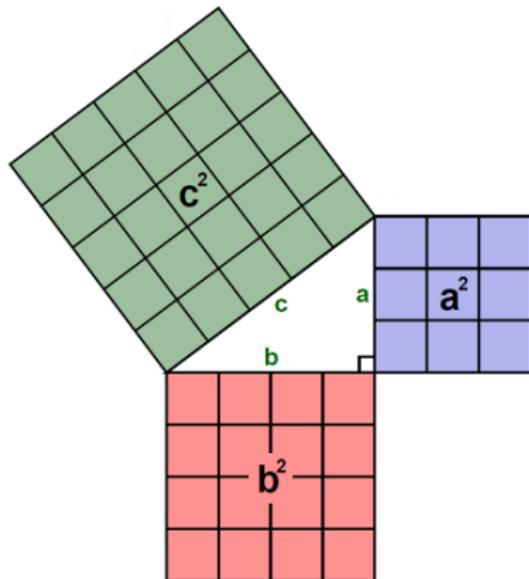
$$\begin{aligned}\text{Var}\left(\frac{1}{\sqrt{2}}X_1 + \frac{1}{\sqrt{2}}X_2\right) &= \left(\frac{1}{\sqrt{2}}\right)^2 \text{Var}(X_1) + \left(\frac{1}{\sqrt{2}}\right)^2 \text{Var}(X_2) \\ &\quad + 2 \frac{1}{\sqrt{2}} \frac{1}{\sqrt{2}} \text{Cov}(X_1, X_2)\end{aligned}$$

$$\begin{aligned}\text{Var}\left(\frac{1}{\sqrt{2}}X_1 + \frac{-1}{\sqrt{2}}X_2\right) &= \left(\frac{1}{\sqrt{2}}\right)^2 \text{Var}(X_1) + \left(\frac{-1}{\sqrt{2}}\right)^2 \text{Var}(X_2) \\ &\quad + 2 \frac{1}{\sqrt{2}} \frac{-1}{\sqrt{2}} \text{Cov}(X_1, X_2)\end{aligned}$$

'Directional', Non
P.D., & Singular

$$\begin{bmatrix} 1 & .6 & .4 \\ .6 & 1 & .4 \\ .4 & .4 & 1 \end{bmatrix} \quad \begin{bmatrix} 1 & 1 & 1 \\ 1 & 1 & .5 \\ 1 & .5 & 1 \end{bmatrix} \quad \begin{bmatrix} 1 & 1 & .5 \\ 1 & 1 & .5 \\ .5 & .5 & 1 \end{bmatrix}$$

PCA as an *approximate* Pythagorean Theorem



$$c^2 = 5^2 + 6^2$$

$$c^2 = c^2 + 3^2$$

$$c^2 = 5^2 + 6^2 + 3^2$$