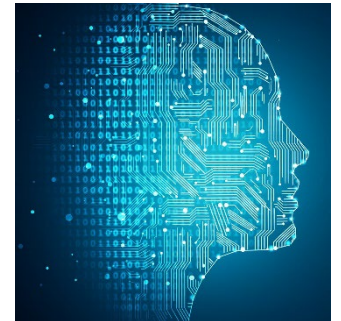


Machine Learning

Principal Components Analysis



Kevin Moon (kevin.moon@usu.edu)
STAT/CS 5810/6655





- Dimensionality Reduction
- Projections
- PCA
 - Projection Perspective
 - Maximum Variance Perspective
 - Connection to SVD
- MDS
 - Connection to PCA

Unsupervised Learning



- On to unsupervised learning!
- We'll come back to supervised learning later.
- There are three main unsupervised learning topics that we will cover in this course:
 1. Dimensionality reduction
 2. Clustering
 3. Density estimation

Dimensionality Reduction



- In dimensionality reduction problems, we observe $\mathbf{x}_1, \dots, \mathbf{x}_n \in \mathbb{R}^d$
 - Notice we do not have labels (unsupervised learning)
- Goal: transform these variables to new ones

$$\mathbf{x}_i \rightarrow \boldsymbol{\theta}_i \in \mathbb{R}^k$$

where $k < d$, such that information loss is minimized

Dimensionality Reduction



Reasons for doing dimensionality reduction:

- Computational efficiency
- Visualization ($k = 2, 3$)
- Compression
- Interpreting data (which dimensions are important?)
- Eliminate rank deficiency
- Eliminate useless/noisy features
- Avoid overfitting

Dimensionality Reduction Types



- Methods for dimensionality reduction can be classified according to:
 1. How is information loss quantified?
 2. Supervised or unsupervised?
 3. Feature selection (select existing features) or feature extraction (form new features from old ones)?
 4. Parametric or nonparametric?
 5. Linear or nonlinear?
 6. Generative or discriminative?



- The first method for dimensionality reduction that we will discuss is principal components analysis (PCA)
 1. How is information loss quantified? **Least squares**
 2. Supervised or **unsupervised**?
 3. Feature selection or **feature extraction**?
 4. **Parametric** or nonparametric?
 5. **Linear** or nonlinear?
 6. Generative or discriminative?

Either, although our initial perspective is discriminative

Importance of PCA



PCA is perhaps the most important unsupervised method

- For high-dimensional supervised learning problems, performing PCA first can be very helpful
- Many other dimensionality reduction methods can be framed as learning a different representation of the data and then applying PCA
 - E.g. diffusion maps and kernel PCA (kernel trick applied to PCA)
- Other dimensionality reduction methods use PCA as a preprocessing step
- You should become very familiar with it



Projections



- Let $\mathbf{a}_1, \dots, \mathbf{a}_k \in \mathbb{R}^d$ be linearly independent column vectors. Denote

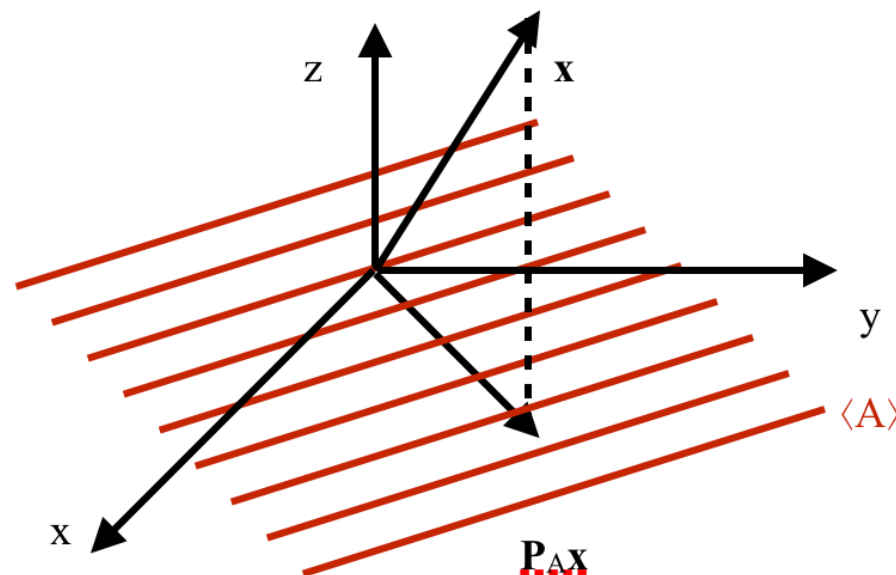
$$\mathbf{A} = [\mathbf{a}_1 \cdots \mathbf{a}_k] \quad (d \times k)$$

- The linear span of $\mathbf{a}_1, \dots, \mathbf{a}_k$ is the column span of \mathbf{A} , written

$$\begin{aligned} \langle \mathbf{A} \rangle &= \text{colspan}(\mathbf{A}) \\ &= \{ \mathbf{A}\boldsymbol{\theta} \mid \boldsymbol{\theta} \in \mathbb{R}^k \}. \end{aligned}$$

- The *projection* onto $\langle \mathbf{A} \rangle$ is the mapping $\mathbf{P}_A : \mathbb{R}^d \rightarrow \langle \mathbf{A} \rangle \subseteq \mathbb{R}^d$

$\mathbf{P}_A \mathbf{x}$ = closest point in $\langle \mathbf{A} \rangle$ to \mathbf{x}



Projections



- Every point in $\langle A \rangle$ equals $A\theta$ for some $\theta \in \mathbb{R}^k$
- Therefore, $P_A x = A\hat{\theta}$, where

$$\hat{\theta} = \arg \min_{\theta} \|x - A\theta\|$$

- We have previously seen that the solution is

$$\hat{\theta} = (A^T A)^{-1} A^T x$$

- Therefore,

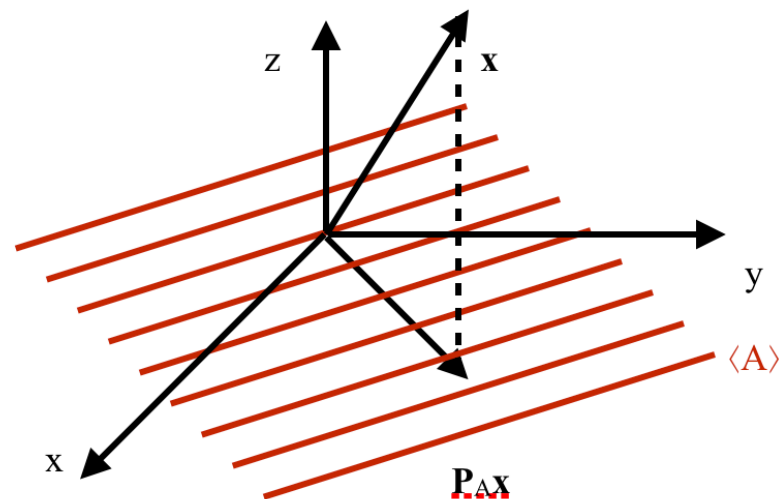
$$\begin{aligned} P_A x &= A\hat{\theta} \\ &= A(A^T A)^{-1} A^T x \\ \Rightarrow P_A &= A(A^T A)^{-1} A^T \end{aligned}$$

Properties of Projections



- If $\mathbf{a}_1, \dots, \mathbf{a}_k$ are orthonormal, then $\mathbf{P}_A = \mathbf{A}\mathbf{A}^T$ ($\mathbf{A}^T\mathbf{A} = \mathbf{I}$)
- The orthogonality principle states that

$$\forall \mathbf{x}, \mathbf{x} - \mathbf{P}_A \mathbf{x} \in \langle \mathbf{A} \rangle^\perp$$



- Proof: Let $\mathbf{u} \in \langle \mathbf{A} \rangle$. Then we can write $\mathbf{u} = \mathbf{A}\boldsymbol{\theta}$ for some $\boldsymbol{\theta} \in \mathbb{R}^k$. Then

$$\begin{aligned} \langle \mathbf{u}, \mathbf{x} - \mathbf{P}_A \mathbf{x} \rangle &= \left\langle \mathbf{u}, \left(\mathbf{I} - \mathbf{A}(\mathbf{A}^T \mathbf{A})^{-1} \mathbf{A}^T \right) \mathbf{x} \right\rangle \\ &= \boldsymbol{\theta}^T \mathbf{A}^T \left(\mathbf{I} - \mathbf{A}(\mathbf{A}^T \mathbf{A})^{-1} \mathbf{A}^T \right) \mathbf{x} \\ &= \boldsymbol{\theta}^T \mathbf{A}^T \mathbf{x} - \boldsymbol{\theta}^T \mathbf{A}^T \mathbf{x} \\ &= 0. \end{aligned}$$

Group Exercise



1. Show that projection matrices are idempotent, that is,

$$\mathbf{P}_A^2 = \mathbf{P}_A.$$

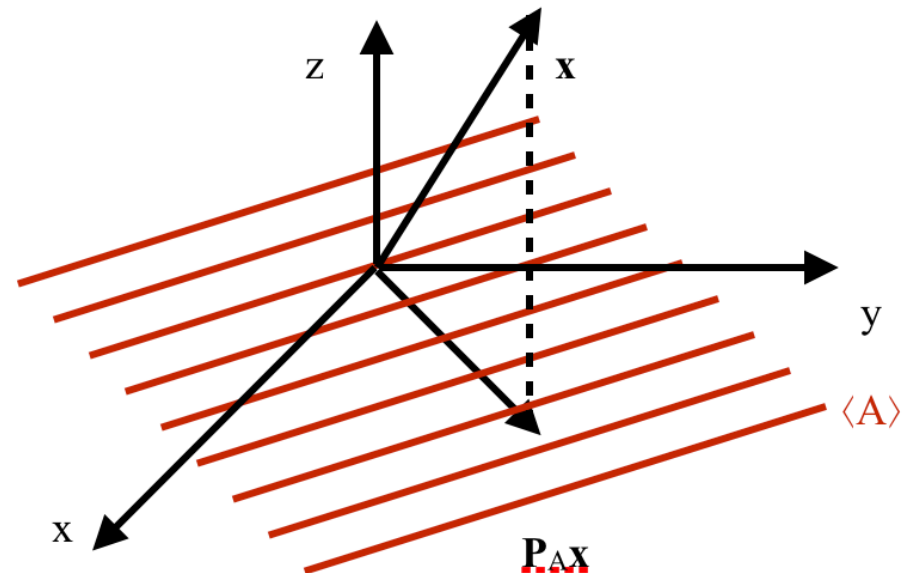
Give an intuitive explanation of this property.

2. Let \mathbf{B} be a $d \times (d - k)$ full rank matrix such that

$$\langle \mathbf{B} \rangle = \langle \mathbf{A} \rangle^\perp.$$

Determine formulas for

- (a) $\mathbf{P}_A \mathbf{P}_B$
- (b) $\mathbf{P}_A + \mathbf{P}_B$





- The idea behind PCA is to approximate

$$\mathbf{x}_i \approx \boldsymbol{\mu} + \mathbf{A}\boldsymbol{\theta}_i$$

where

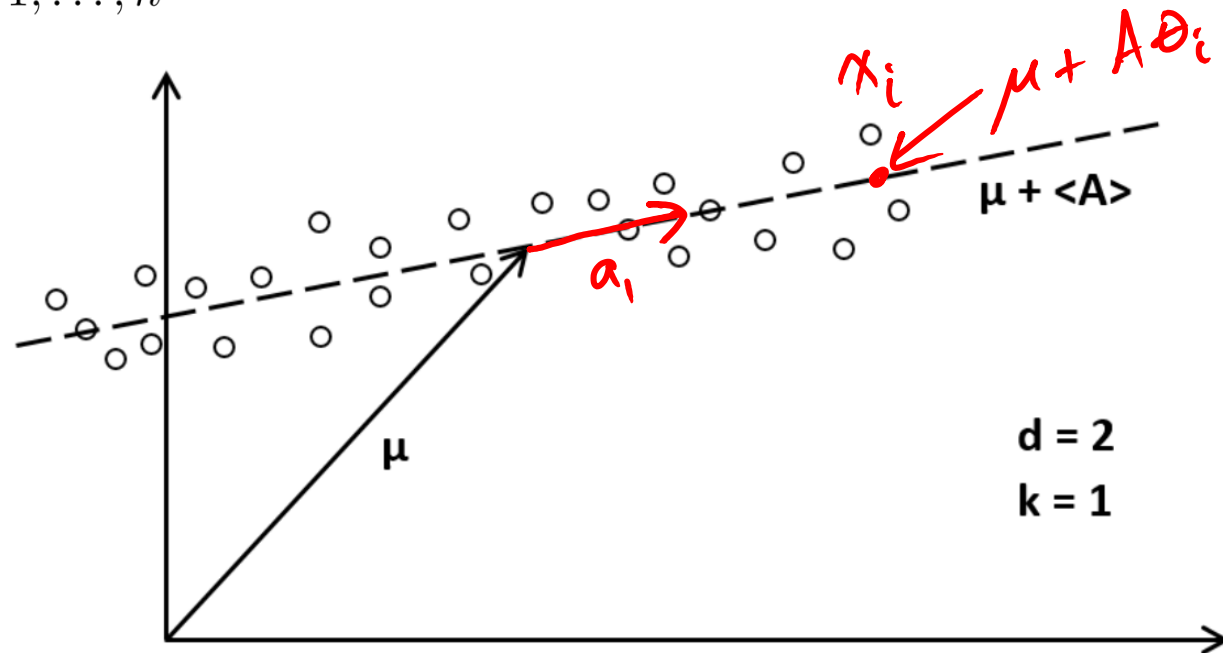
$$\boldsymbol{\mu} \in \mathbb{R}^d$$

$$\mathbf{A} \in \mathcal{A}_k := \{\mathbf{A} \in \mathbb{R}^{d \times k} \mid \mathbf{A}^T \mathbf{A} = \mathbf{I}_{k \times k}\}$$

$$\boldsymbol{\theta}_i \in \mathbb{R}^k, i = 1, \dots, n$$

$$\mathbf{A} = [\mathbf{a}_1] \in \mathbb{R}^d$$

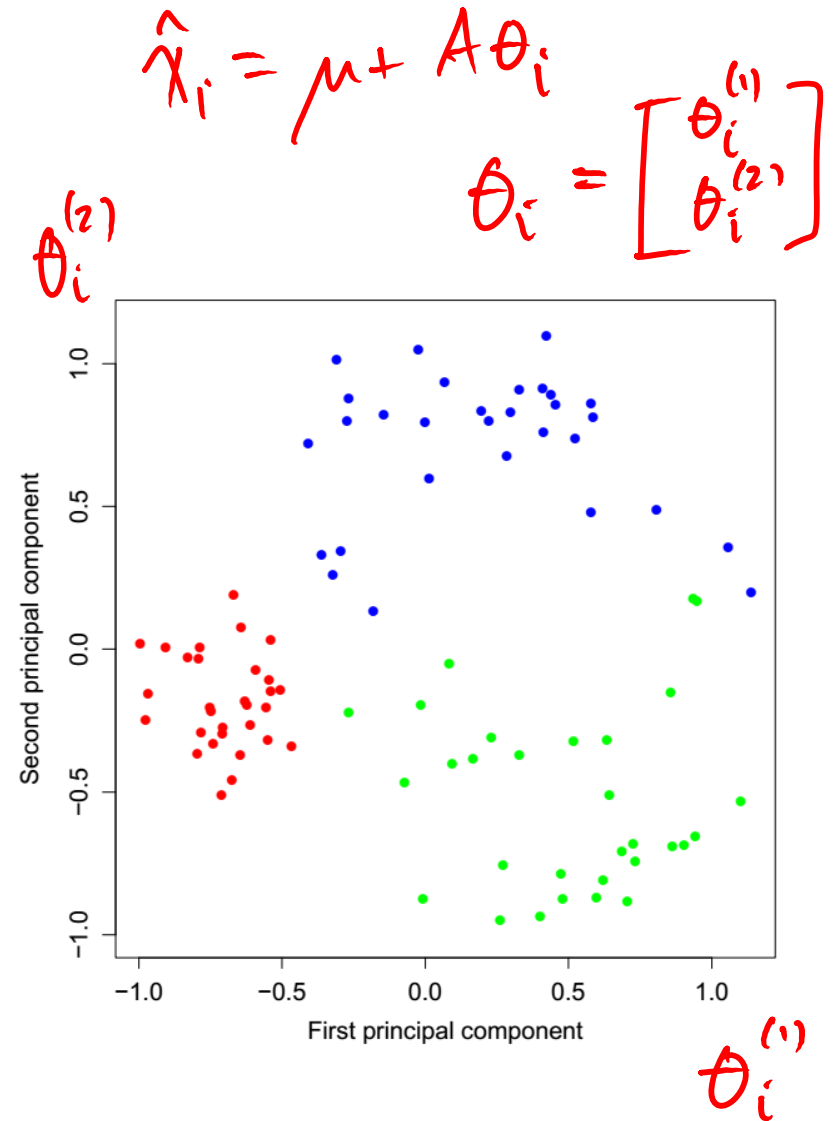
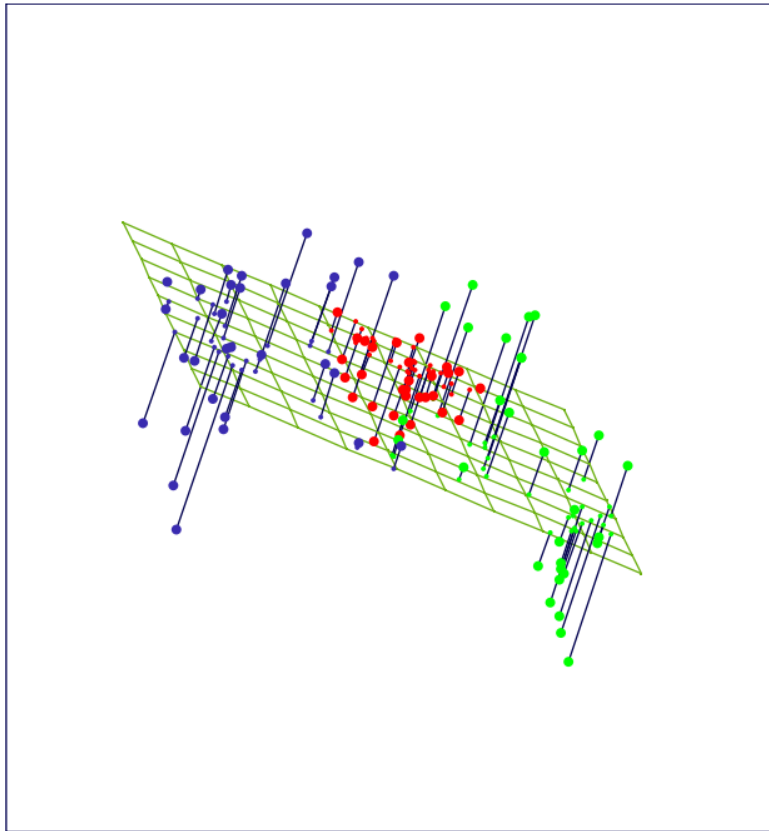
- Example:** $d = 2, k = 1$.



PCA



- **Example:** $d = 3, k = 2$





- Mathematically, we define $\boldsymbol{\mu}, \mathbf{A}, \boldsymbol{\theta}_1, \dots, \boldsymbol{\theta}_n$ to be the solution of

$$\min_{\substack{\boldsymbol{\mu}, \mathbf{A}, \boldsymbol{\theta}_1, \dots, \boldsymbol{\theta}_n \\ \mathbf{A}^T \mathbf{A} = \mathbf{I}}} \sum_{i=1}^n \|\mathbf{x}_i - \boldsymbol{\mu} - \mathbf{A}\boldsymbol{\theta}_i\|^2$$

- PCA gives the least squares rank- k linear approximation to the data set.
- The solution is given in terms of the spectral (or eigenvalue decomposition) of the sample covariance matrix:

$$\mathbf{S} = \frac{1}{n} \sum_{i=1}^n (\mathbf{x}_i - \bar{\mathbf{x}})(\mathbf{x}_i - \bar{\mathbf{x}})^T = \mathbf{U}\boldsymbol{\Lambda}\mathbf{U}^T$$



- Note that S is PSD so $\lambda_i \geq 0$:

$$\mathbf{z}^T S \mathbf{z} = \frac{1}{n} \sum_{i=1}^n \mathbf{z}^T (\mathbf{x}_i - \bar{\mathbf{x}}) (\mathbf{x}_i - \bar{\mathbf{x}})^T \mathbf{z}$$

$$= \frac{1}{n} \sum_{i=1}^n (\mathbf{z}^T (\mathbf{x}_i - \bar{\mathbf{x}}))^2 \geq 0$$

- A solution to PCA is:

$$\begin{aligned} \boldsymbol{\mu} &= \bar{\mathbf{x}}, & \mathbf{A} &= [\mathbf{u}_1, \dots, \mathbf{u}_k] \\ \boldsymbol{\theta}_i &= \mathbf{A}^T (\mathbf{x}_i - \bar{\mathbf{x}}) \end{aligned}$$

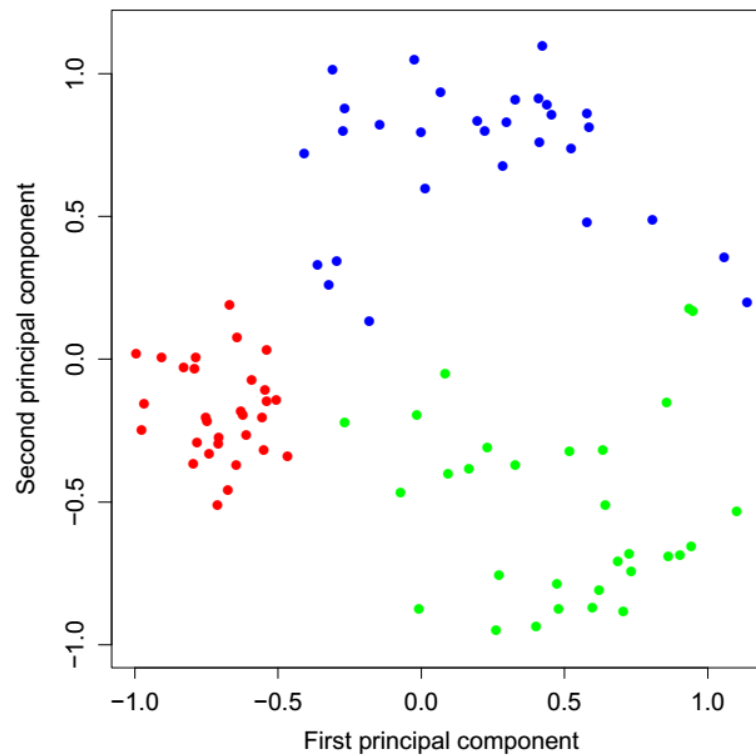
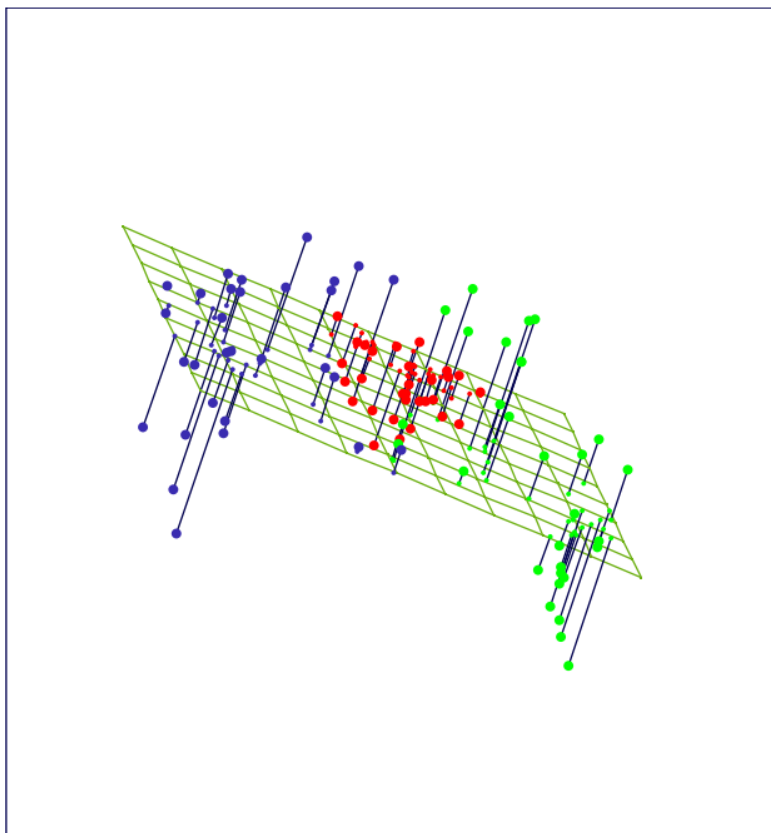
Terminology and Concepts



- Principal components
 - $\boldsymbol{\theta}^{(j)}$ = j th principal component
- Principal eigenvectors/directions
 - \boldsymbol{u}_j = j th principal eigenvector/direction
- Reconstruction of \boldsymbol{x}_i
 - $\hat{\boldsymbol{x}}_i = \boldsymbol{\mu} + \boldsymbol{A}\boldsymbol{\theta}_i$



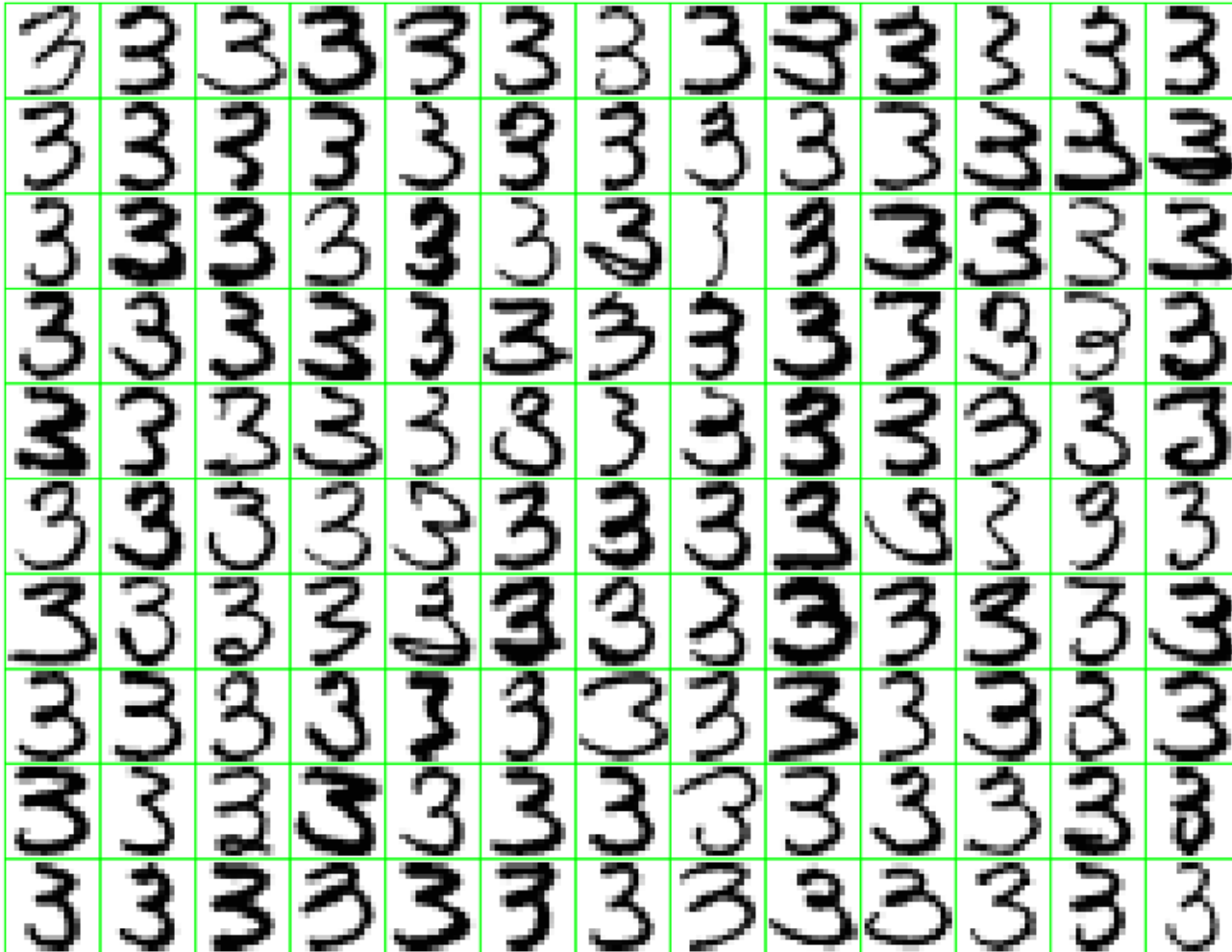
- **Example:** $d = 3, k = 2$



Example: Handwritten Digits



- Training data



Example: Handwritten Digits



- Reconstruction

reconstructed with 2 bases

$k=2$



reconstructed with 10 bases

$k=10$



reconstructed with 100 bases

$k=100$



reconstructed with 506 bases



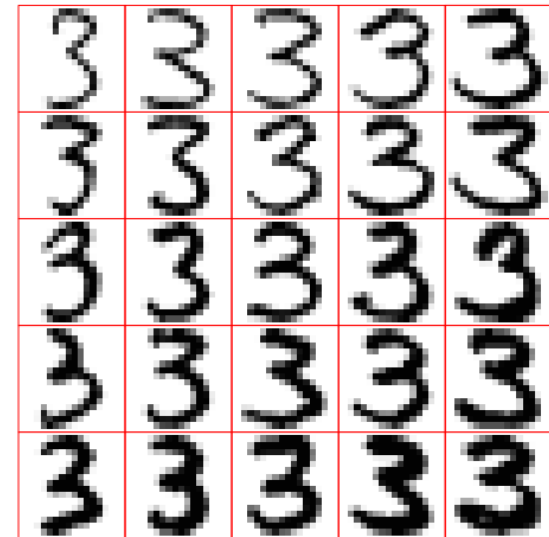
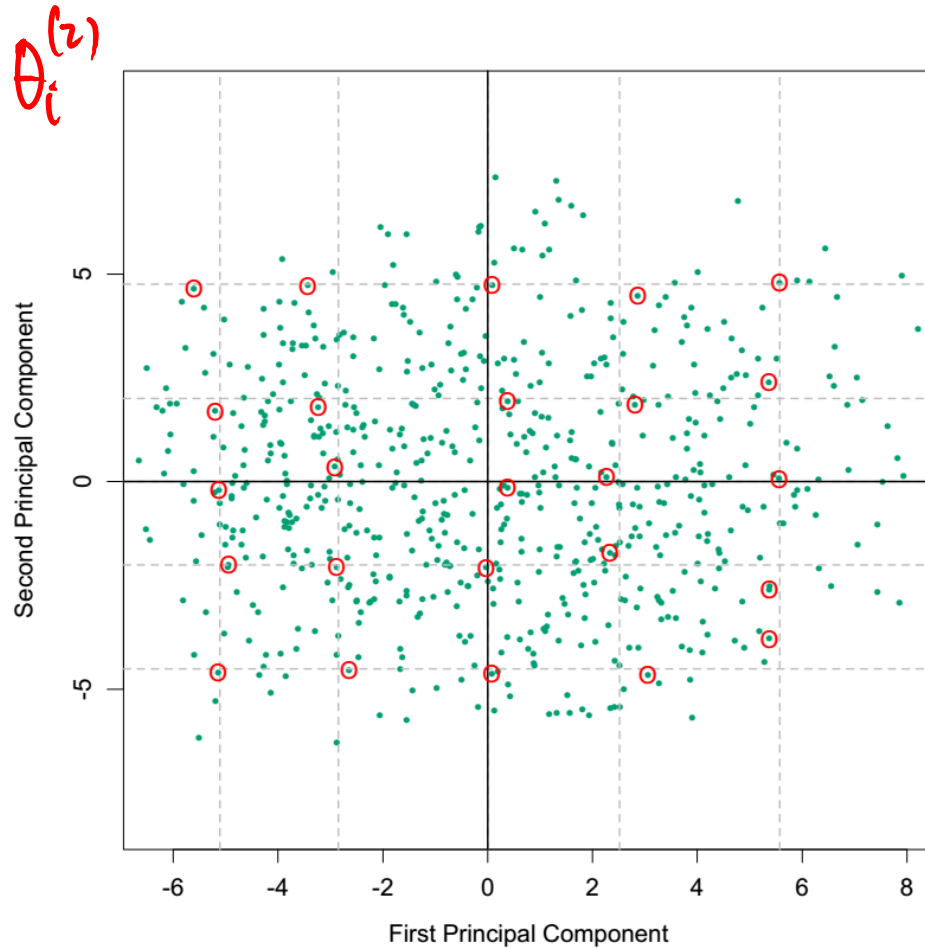
$A\theta$

$$\hat{x}_i = \mu + \theta_i^{(1)} u_1 + \theta_i^{(2)} u_2 = \boxed{3} + \theta_i^{(1)} \boxed{3} + \theta_i^{(2)} \boxed{3}$$

Example: Handwritten Digits



- First two PCs



$\theta_i^{(1)}$

Connection to Projections

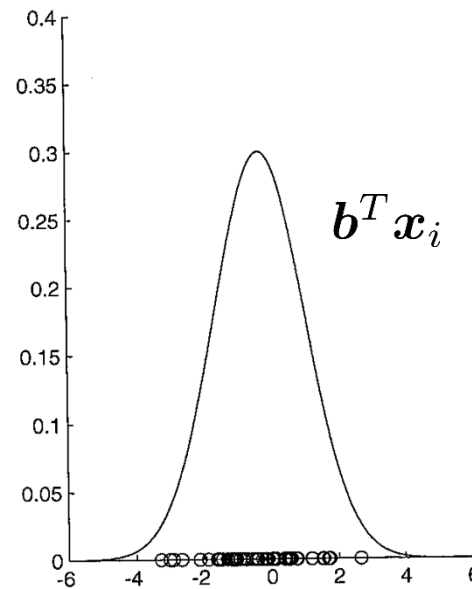
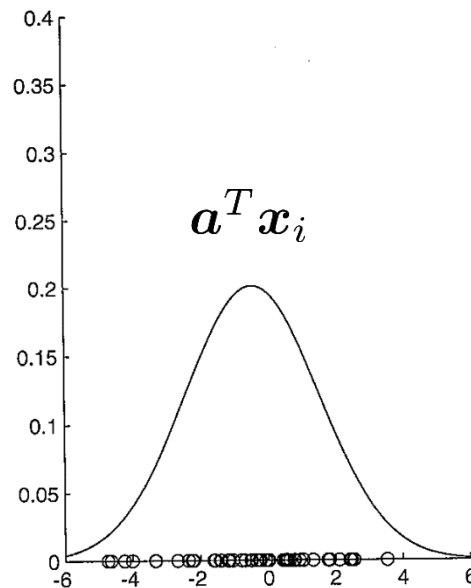
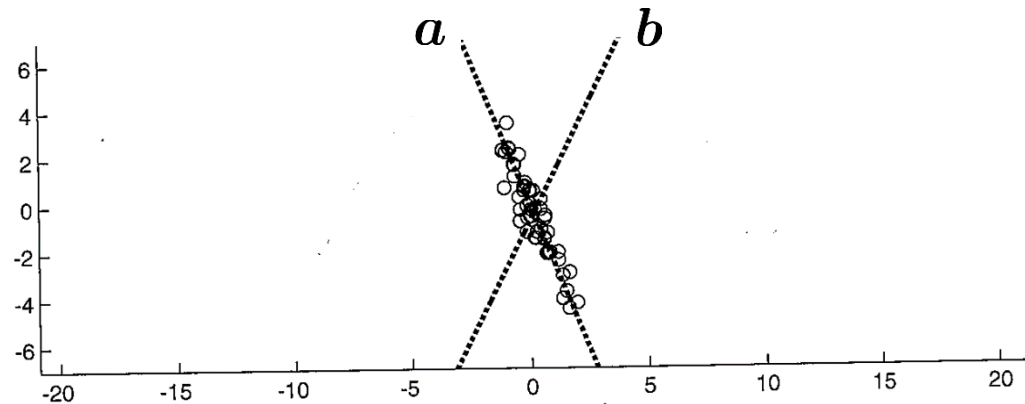


- Suppose $\bar{\mathbf{x}} = \mathbf{0}$
- Recall that due to orthonormality, the projection matrix is $\mathbf{A}\mathbf{A}^T$
- The rank- k approximation to \mathbf{x}_i is

$$\begin{aligned}\hat{\mathbf{x}}_i &= \mathbf{A}\boldsymbol{\theta}_i = \mathbf{A}\mathbf{A}^T \mathbf{x}_i \\ &= \sum_{j=1}^k \mathbf{u}_j \theta_i^{(j)} \\ \theta_i^{(j)} &= \mathbf{u}_j^T \mathbf{x}_i\end{aligned}$$

- Intuition:
 - Columns of \mathbf{A} define a k dimensional coordinate system for $\langle \mathbf{A} \rangle$
 - $\boldsymbol{\theta}_i = \mathbf{A}^T \mathbf{x}_i$ are the coordinates of $\hat{\mathbf{x}}_i$ in the subspace

Connection to Projections



Maximum Variance Perspective



- Suppose $\bar{\mathbf{x}} = \mathbf{0}$
- Let \mathbf{X} be a random vector of which $\mathbf{x}_1, \dots, \mathbf{x}_n$ are realizations
- Goal: find the unit vector $\mathbf{a}_1 \in \mathbb{R}^d$ ($\|\mathbf{a}_1\| = 1$) which maximizes the sample variance of

$$\theta^{(1)} = \mathbf{a}_1^T \mathbf{X}$$

- Sample mean of $\theta^{(1)}$:

$$\frac{1}{n} \sum_{i=1}^n \mathbf{a}_1^T \mathbf{x}_i = \mathbf{a}_1^T \left(\frac{1}{n} \sum_{i=1}^n \mathbf{x}_i \right) = \mathbf{a}_1^T \bar{\mathbf{x}}$$



Maximum Variance Perspective



- Sample variance of $\theta^{(1)}$:

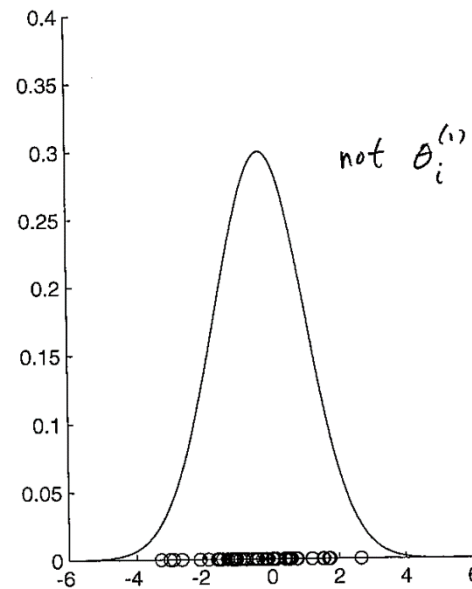
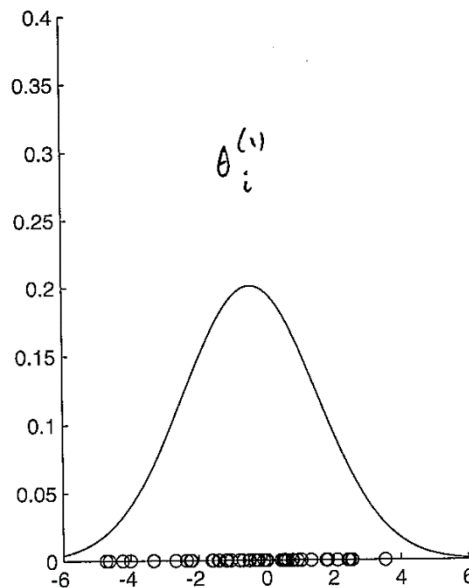
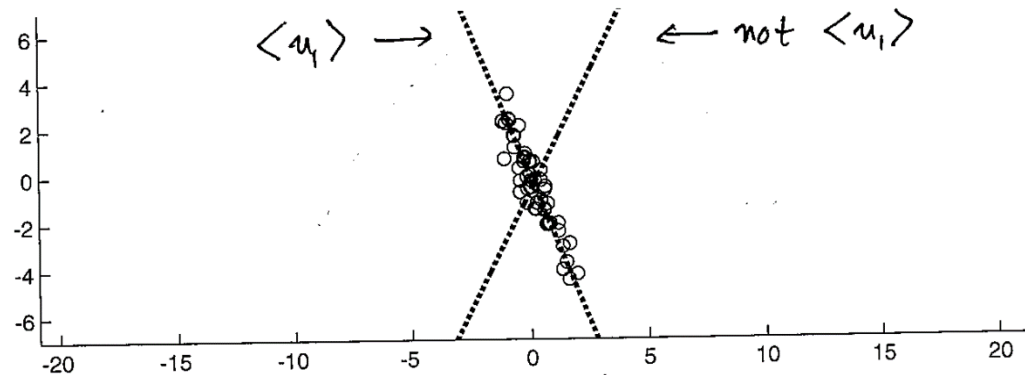
$$\begin{aligned}\frac{1}{n} \sum_{i=1}^n (\mathbf{a}_1^T \mathbf{x}_i)^2 &= \frac{1}{n} \sum_{i=1}^n (\mathbf{a}_1^T \mathbf{x}_i)(\mathbf{x}_i^T \mathbf{a}_1) \\ &= \mathbf{a}_1^T \left(\frac{1}{n} \sum_{i=1}^n \mathbf{x}_i \mathbf{x}_i^T \right) \mathbf{a}_1 \\ &= \mathbf{a}_1^T S \mathbf{a}_1\end{aligned}$$

- The solution of

$$\max_{\mathbf{a}_1: \|\mathbf{a}_1\|=1} \mathbf{a}_1^T S \mathbf{a}_1$$

is \mathbf{u}_1 (the first eigenvector of S). Thus $\theta^{(1)} = \mathbf{u}_1^T \mathbf{X}$.

Maximum Variance Perspective

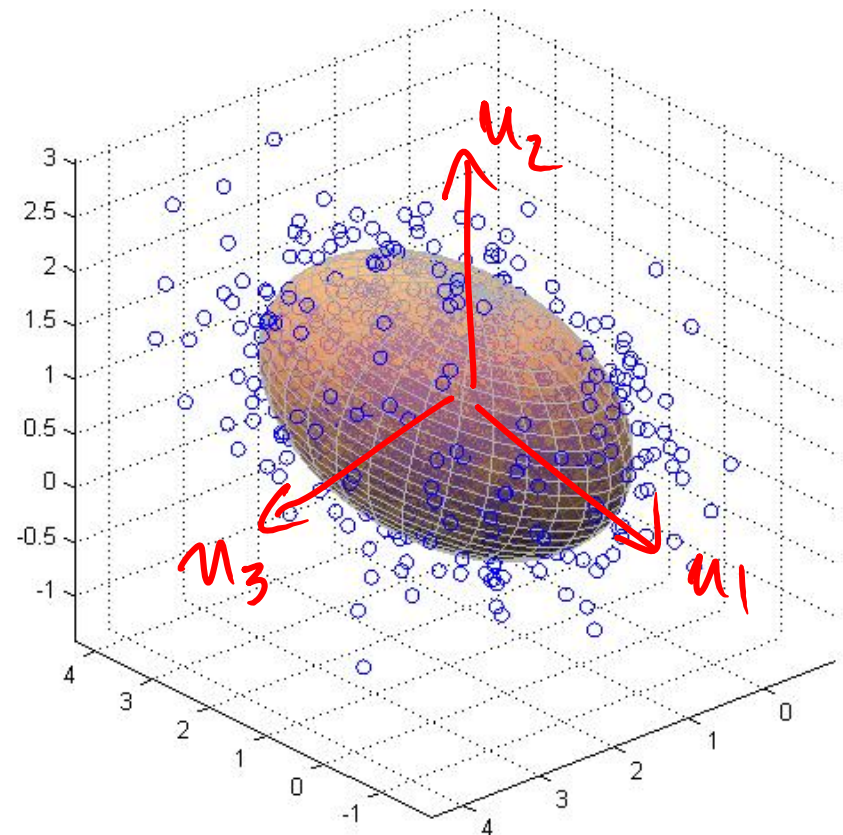




Maximum Variance Perspective



- More generally, we have the following result:
- **Theorem:** Let $\theta^{(k)} = \mathbf{a}_k^T \mathbf{X}$ and $\text{var}(\theta^{(k)}) = \frac{1}{n} \sum_{i=1}^n (\mathbf{a}_k^T \mathbf{x}_i)^2$. A vector \mathbf{a}_k that maximizes $\text{var}(\theta^{(k)})$ subject to
 - $\|\mathbf{a}_k\| = 1$
 - $\mathbf{a}_k \perp \mathbf{u}_1, \dots, \mathbf{u}_{k-1}$is $\mathbf{a}_k = \mathbf{u}_k$.
- What is the optimized variance of $\theta^{(k)}$?



Optimized Variance



$$\text{var}_{\text{samp}}(\mathbf{u}_j^T \mathbf{X}) = \mathbf{u}_j^T \mathbf{S} \mathbf{u}_j$$

$$= \mathbf{u}_j^T \mathbf{U} \mathbf{\Lambda} \mathbf{U}^T \mathbf{u}_j$$

$$= [0, \dots, 0, 1, 0, \dots, 0] \mathbf{\Lambda} \begin{bmatrix} 0 \\ \vdots \\ 0 \\ 1 \\ 0 \\ \vdots \\ 0 \end{bmatrix}$$

$$= \lambda_j$$

Selecting k



- It can be shown that the optimal objective function for PCA is

$$\min_{\substack{\mu, A, \theta_1, \dots, \theta_n \\ A^T A = I}} \sum_{i=1}^n \|x_i - \mu - A\theta_i\|^2 = n(\lambda_{k+1} + \dots + \lambda_d)$$

- When $k = 0$, this specializes to

$$\min_{\mu} \sum_{i=1}^n \|x_i - \mu\|^2 = n(\lambda_1 + \dots + \lambda_d)$$

which we call the total variation of the data.

- Common heuristic: choose smallest k s.t.

$$\frac{\lambda_1 + \dots + \lambda_k}{\lambda_1 + \dots + \lambda_d} = \% \text{ of variance explained by the top } k \text{ PCs}$$

- Common threshold is 95%

Connection to SVD



- Assume $\bar{\mathbf{x}} = \mathbf{0}$
- Data matrix ($d \times n$) $X = [\mathbf{x}_1, \dots, \mathbf{x}_n]$
 - Recall that $S = \frac{1}{n}XX^T$
- The singular value decomposition (SVD) of X is

$$X = U\Sigma V^T$$

where U ($d \times d$) and V ($n \times n$) are orthogonal matrices, and $\Sigma = \text{diag}(\sigma_1, \dots, \sigma_{\min\{d,n\}})$ is $d \times n$

- Then $u_j = j\text{th left singular vector} = j\text{th principal component}$
 - Also, $\lambda_j = \frac{1}{n}\sigma_j^2$

Multidimensional Scaling (MDS)

Multidimensional Scaling (MDS)



- Another dimensionality reduction method, typically focused on visualization
- **Example:** Suppose we have the following distance matrix between major cities in the US (3D distances)

		1	2	3	4	5	6	7	8	9
		BOST	NY	DC	MIAM	CHIC	SEAT	SF	LA	DENV
1	BOSTON	0	206	429	1504	963	2976	3095	2979	1949
2	NY	206	0	233	1308	802	2815	2934	2786	1771
3	DC	429	233	0	1075	671	2684	2799	2631	1616
4	MIAMI	1504	1308	1075	0	1329	3273	3053	2687	2037
5	CHICAGO	963	802	671	1329	0	2013	2142	2054	996
6	SEATTLE	2976	2815	2684	3273	2013	0	808	1131	1307
7	SF	3095	2934	2799	3053	2142	808	0	379	1235
8	LA	2979	2786	2631	2687	2054	1131	379	0	1059
9	DENVER	1949	1771	1616	2037	996	1307	1235	1059	0

- Can we create a 2D representation of this dataset?
- Use MDS

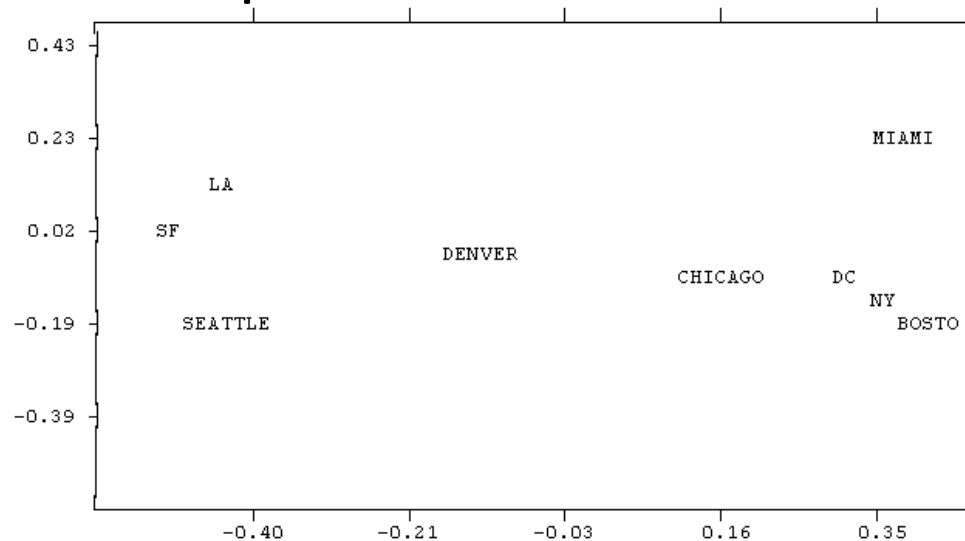
MDS Example



- Distance matrix

	1	2	3	4	5	6	7	8	9
	BOST	NY	DC	MIAM	CHIC	SEAT	SF	LA	DENV
1	0	206	429	1504	963	2976	3095	2979	1949
2	206	0	233	1308	802	2815	2934	2786	1771
3	429	233	0	1075	671	2684	2799	2631	1616
4	1504	1308	1075	0	1329	3273	3053	2687	2037
5	963	802	671	1329	0	2013	2142	2054	996
6	2976	2815	2684	3273	2013	0	808	1131	1307
7	3095	2934	2799	3053	2142	808	0	379	1235
8	2979	2786	2631	2687	2054	1131	379	0	1059
9	1949	1771	1616	2037	996	1307	1235	1059	0

- MDS-generated map





- MDS works by minimizing a loss function between the distances in the high-dimensional space and the distances in the low-dimensional space
- Setup: $\mathbf{x}_1, \dots, \mathbf{x}_n \in \mathbb{R}^d$ are the measurements in the original, high –dimensional space
- Find representations $\mathbf{y}_1, \dots, \mathbf{y}_n \in \mathbb{R}^k$ ($k < d$) such that the loss function is minimized
- Different loss functions give you different variations on MDS



- $d_{ij} = \|\mathbf{x}_i - \mathbf{x}_j\|$ (pairwise Euclidean distance)
- Metric MDS minimizes the “Stress” function:

$$\text{Stress}_M(\mathbf{y}_1, \dots, \mathbf{y}_n) = \left(\frac{\sum_{i=1}^n \sum_{j=1}^n (d_{ij} - \|\mathbf{y}_i - \mathbf{y}_j\|)^2}{\sum_{i=1}^n \sum_{j=1}^n d_{ij}^2} \right)^{1/2}$$

- The \mathbf{y}_i s are the variables that are chosen to minimize the stress.
- Other kinds of distances can be chosen to give different variations

Choosing the dimension k

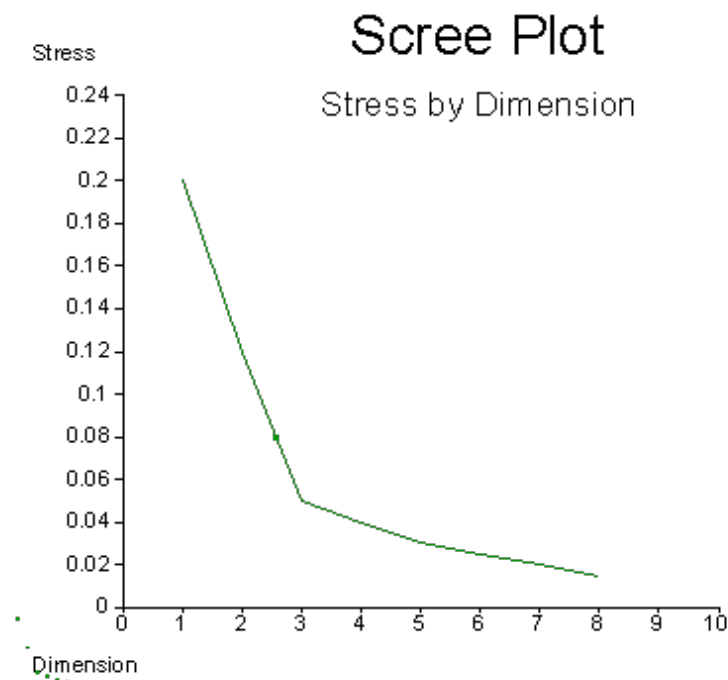


- If the stress is zero, then there is no distortion of the distances
 - Thus smaller stress \Rightarrow a better representation
 - Stress decreases as the dimension k increases (less distortion required with higher dimensions)
- However, if the data are noisy, some “distortion” may be ok
- **Example:** distances from buildings in NYC measured from center of the roof
 - Clearly a 3D dataset
 - But a 3D MDS representation may have nonzero stress/loss

Choosing the dimension k



- How do we know what k should be?
- No surefire answer...
- Common approach is to look at a “scree plot”
- “Elbows” aren’t always obvious and so other approaches may be needed
 - E.g., Shepard diagrams (plot of x_i distances vs y_i distances) may be useful



Nonmetric MDS



- Recall the stress function for metric MDS:

$$\text{Stress}_M(\mathbf{y}_1, \dots, \mathbf{y}_n) = \left(\frac{\sum_{i=1}^n \sum_{j=1}^n (d_{ij} - \|\mathbf{y}_i - \mathbf{y}_j\|)^2}{\sum_{i=1}^n \sum_{j=1}^n d_{ij}^2} \right)^{1/2}$$

- Now let d_{ij} be a measure of dissimilarity (could be Euclidean distance) between the points \mathbf{x}_i and \mathbf{x}_j
- Let f be some monotonic function
- Nonmetric MDS minimizes the following:

$$\begin{aligned} & \text{Stress}_{NM}(f, \mathbf{y}_1, \dots, \mathbf{y}_n) \\ &= \left(\frac{\sum_{i=1}^n \sum_{j=1}^n (f(d_{ij}) - \|\mathbf{y}_i - \mathbf{y}_j\|)^2}{\sum_{i=1}^n \sum_{j=1}^n \|\mathbf{y}_i - \mathbf{y}_j\|^2} \right)^{1/2} \end{aligned}$$

Classical MDS



- A special case of metric MDS
- Let $d_{ij} = \|\mathbf{x}_i - \mathbf{x}_j\|$
- Let $D^{(2)} = [d_{ij}^2]$ be the matrix of squared distances
- Double center the matrix: $B = -\frac{1}{2}JD^{(2)}J$
 - $J = I - \frac{1}{n}\mathbf{1}\mathbf{1}^T$ ($\mathbf{1}$ is an n -dimensional vector of ones)
- Define b_{ij} to be the i, j th entry of B
- Classical MDS minimizes the “Strain” function:
$$Strain(y_1, \dots, y_n) = \left(\frac{\sum_{i,j} (b_{ij} - \langle y_i, y_j \rangle)^2}{\sum_{i,j} b_{ij}^2} \right)^{1/2}$$

Classical MDS



- Classical MDS minimizes the “Strain” function:

$$\text{Strain}(y_1, \dots, y_n) = \left(\frac{\sum_{i,j} (b_{ij} - \langle y_i, y_j \rangle)^2}{\sum_{i,j} b_{ij}^2} \right)^{1/2}$$

- It turns out, this can be solved using eigenvalue decomposition
 - Efficient computation!
 - Strictly assumes that the distances are Euclidean (may be too restrictive sometimes)
- Can be shown that classical MDS and PCA are equivalent (recover one from the other)

Further Reading



- ISL Section 10.2
- ESL Sections 3.4.1, 14.5, and 14.8

- MDS figures:

<http://www.analytictech.com/borgatti/mds.htm>