

# Text as Data

Justin Grimmer

Associate Professor  
Department of Political Science  
Stanford University

October 14th, 2014

# Finding Lower Dimensional Embeddings of Text

1) Task:

# Finding Lower Dimensional Embeddings of Text

## 1) Task:

- Embed our documents in a lower dimensional space

# Finding Lower Dimensional Embeddings of Text

## 1) Task:

- Embed our documents in a lower dimensional space
- Visualize our documents

# Finding Lower Dimensional Embeddings of Text

## 1) Task:

- Embed our documents in a lower dimensional space
- Visualize our documents
- Inference about similarity

# Finding Lower Dimensional Embeddings of Text

## 1) Task:

- Embed our documents in a lower dimensional space
- Visualize our documents
- Inference about similarity
- **Social science inference about behavior**

# Finding Lower Dimensional Embeddings of Text

## 1) Task:

- Embed our documents in a lower dimensional space
- Visualize our documents
- Inference about similarity
- **Social science inference about behavior**

## 2) Objective Function, $f(\mathbf{X}, \theta)$

# Finding Lower Dimensional Embeddings of Text

## 1) Task:

- Embed our documents in a lower dimensional space
- Visualize our documents
- Inference about similarity
- **Social science inference about behavior**

## 2) Objective Function, $f(\mathbf{X}, \theta)$

- Reconstruction error: measure the distance between the **embedded** points and the original points



# Finding Lower Dimensional Embeddings of Text

## 1) Task:

- Embed our documents in a lower dimensional space
- Visualize our documents
- Inference about similarity
- **Social science inference about behavior**

## 2) Objective Function, $f(\mathbf{X}, \theta)$

- Reconstruction error: measure the distance between the **embedded** points and the original points  $\rightsquigarrow \theta$  will contain basis vectors (**principal components**) and **weights** to those vectors

# Finding Lower Dimensional Embeddings of Text

## 1) Task:

- Embed our documents in a lower dimensional space
- Visualize our documents
- Inference about similarity
- **Social science inference about behavior**

## 2) Objective Function, $f(\mathbf{X}, \theta)$

- Reconstruction error: measure the distance between the **embedded** points and the original points  $\rightsquigarrow \theta$  will contain basis vectors (**principal components**) and **weights** to those vectors
- Distance preservation: measure how well pairwise distances are preserved with distances  $\rightsquigarrow \theta$  a set of new lower-dimensional coordinates

# Finding Lower Dimensional Embeddings of Text

## 1) Task:

- Embed our documents in a lower dimensional space
- Visualize our documents
- Inference about similarity
- **Social science inference about behavior**

## 2) Objective Function, $f(\mathbf{X}, \theta)$

- Reconstruction error: measure the distance between the **embedded** points and the original points  $\rightsquigarrow \theta$  will contain basis vectors (**principal components**) and **weights** to those vectors
- Distance preservation: measure how well pairwise distances are preserved with distances  $\rightsquigarrow \theta$  a set of new lower-dimensional coordinates

## 3) Optimization

# Finding Lower Dimensional Embeddings of Text

## 1) Task:

- Embed our documents in a lower dimensional space
- Visualize our documents
- Inference about similarity
- **Social science inference about behavior**

## 2) Objective Function, $f(\mathbf{X}, \theta)$

- Reconstruction error: measure the distance between the **embedded** points and the original points  $\rightsquigarrow \theta$  will contain basis vectors (**principal components**) and **weights** to those vectors
- Distance preservation: measure how well pairwise distances are preserved with distances  $\rightsquigarrow \theta$  a set of new lower-dimensional coordinates

## 3) Optimization $\rightsquigarrow$ eigenvectors, eigenvalues

# Finding Lower Dimensional Embeddings of Text

## 1) Task:

- Embed our documents in a lower dimensional space
- Visualize our documents
- Inference about similarity
- **Social science inference about behavior**

## 2) Objective Function, $f(\mathbf{X}, \theta)$

- Reconstruction error: measure the distance between the **embedded** points and the original points  $\rightsquigarrow \theta$  will contain basis vectors (**principal components**) and **weights** to those vectors
- Distance preservation: measure how well pairwise distances are preserved with distances  $\rightsquigarrow \theta$  a set of new lower-dimensional coordinates

## 3) Optimization $\rightsquigarrow$ eigenvectors, eigenvalues

- Reconstruction error  $\rightsquigarrow$  Search over components and weights

# Finding Lower Dimensional Embeddings of Text

## 1) Task:

- Embed our documents in a lower dimensional space
- Visualize our documents
- Inference about similarity
- **Social science inference about behavior**

## 2) Objective Function, $f(\mathbf{X}, \theta)$

- Reconstruction error: measure the distance between the **embedded** points and the original points  $\rightsquigarrow \theta$  will contain basis vectors (**principal components**) and **weights** to those vectors
- Distance preservation: measure how well pairwise distances are preserved with distances  $\rightsquigarrow \theta$  a set of new lower-dimensional coordinates

## 3) Optimization $\rightsquigarrow$ eigenvectors, eigenvalues

- Reconstruction error  $\rightsquigarrow$  Search over components and weights
- Distance preservation  $\rightsquigarrow$  Search over locations

# Finding Lower Dimensional Embeddings of Text

## 1) Task:

- Embed our documents in a lower dimensional space
- Visualize our documents
- Inference about similarity
- **Social science inference about behavior**

## 2) Objective Function, $f(\mathbf{X}, \theta)$

- Reconstruction error: measure the distance between the **embedded** points and the original points  $\rightsquigarrow \theta$  will contain basis vectors (**principal components**) and **weights** to those vectors
- Distance preservation: measure how well pairwise distances are preserved with distances  $\rightsquigarrow \theta$  a set of new lower-dimensional coordinates

## 3) Optimization $\rightsquigarrow$ eigenvectors, eigenvalues

- Reconstruction error  $\rightsquigarrow$  Search over components and weights
- Distance preservation  $\rightsquigarrow$  Search over locations

## 4) Validation

# Finding Lower Dimensional Embeddings of Text

## 1) Task:

- Embed our documents in a lower dimensional space
- Visualize our documents
- Inference about similarity
- **Social science inference about behavior**

## 2) Objective Function, $f(\mathbf{X}, \theta)$

- Reconstruction error: measure the distance between the **embedded** points and the original points  $\rightsquigarrow \theta$  will contain basis vectors (**principal components**) and **weights** to those vectors
- Distance preservation: measure how well pairwise distances are preserved with distances  $\rightsquigarrow \theta$  a set of new lower-dimensional coordinates

## 3) Optimization $\rightsquigarrow$ eigenvectors, eigenvalues

- Reconstruction error  $\rightsquigarrow$  Search over components and weights
- Distance preservation  $\rightsquigarrow$  Search over locations

## 4) Validation

- a) Labeling exercise: What are the dimensions?



# Finding Lower Dimensional Embeddings of Text

## 1) Task:

- Embed our documents in a lower dimensional space
- Visualize our documents
- Inference about similarity
- **Social science inference about behavior**

## 2) Objective Function, $f(\mathbf{X}, \theta)$

- Reconstruction error: measure the distance between the **embedded** points and the original points  $\rightsquigarrow \theta$  will contain basis vectors (**principal components**) and **weights** to those vectors
- Distance preservation: measure how well pairwise distances are preserved with distances  $\rightsquigarrow \theta$  a set of new lower-dimensional coordinates

## 3) Optimization $\rightsquigarrow$ eigenvectors, eigenvalues

- Reconstruction error  $\rightsquigarrow$  Search over components and weights
- Distance preservation  $\rightsquigarrow$  Search over locations

## 4) Validation

- a) Labeling exercise: What are the dimensions?
- b) Are the dimensions interesting semantically?

# An Introduction to Eigenvectors, Values, and Diagonalization

## Definition

Suppose  $\mathbf{A}$  is an  $N \times N$  matrix and  $\lambda$  is a scalar.

If

$$\mathbf{Ax} = \lambda\mathbf{x}$$

Then  $\mathbf{x}$  is an *eigenvector* and  $\lambda$  is the associated *eigenvalue*

# An Introduction to Eigenvectors, Values, and Diagonalization

## Definition

Suppose  $\mathbf{A}$  is an  $N \times N$  matrix and  $\lambda$  is a scalar.

If

$$\mathbf{Ax} = \lambda \mathbf{x}$$

Then  $\mathbf{x}$  is an *eigenvector* and  $\lambda$  is the associated *eigenvalue*

- $\mathbf{A}$  stretches the eigenvector  $\mathbf{x}$

# An Introduction to Eigenvectors, Values, and Diagonalization

## Definition

Suppose  $\mathbf{A}$  is an  $N \times N$  matrix and  $\lambda$  is a scalar.

If

$$\mathbf{Ax} = \lambda\mathbf{x}$$

Then  $\mathbf{x}$  is an *eigenvector* and  $\lambda$  is the associated *eigenvalue*

- $\mathbf{A}$  stretches the eigenvector  $\mathbf{x}$
- $\mathbf{A}$  stretches  $\mathbf{x}$  by  $\lambda$

# An Introduction to Eigenvectors, Values, and Diagonalization

## Definition

Suppose  $\mathbf{A}$  is an  $N \times N$  matrix and  $\lambda$  is a scalar.

If

$$\mathbf{Ax} = \lambda\mathbf{x}$$

Then  $\mathbf{x}$  is an *eigenvector* and  $\lambda$  is the associated *eigenvalue*

- $\mathbf{A}$  stretches the eigenvector  $\mathbf{x}$
- $\mathbf{A}$  stretches  $\mathbf{x}$  by  $\lambda$
- To find eigenvectors/values: (eigen in R)

# An Introduction to Eigenvectors, Values, and Diagonalization

## Definition

Suppose  $\mathbf{A}$  is an  $N \times N$  matrix and  $\lambda$  is a scalar.

If

$$\mathbf{Ax} = \lambda \mathbf{x}$$

Then  $\mathbf{x}$  is an *eigenvector* and  $\lambda$  is the associated *eigenvalue*

- $\mathbf{A}$  stretches the eigenvector  $\mathbf{x}$
- $\mathbf{A}$  stretches  $\mathbf{x}$  by  $\lambda$
- To find eigenvectors/values: (eigen in R)
  - Find  $\lambda$  that solves  $\det(\mathbf{A} - \lambda \mathbf{I}) = 0$

# An Introduction to Eigenvectors, Values, and Diagonalization

## Definition

Suppose  $\mathbf{A}$  is an  $N \times N$  matrix and  $\lambda$  is a scalar.

If

$$\mathbf{Ax} = \lambda\mathbf{x}$$

Then  $\mathbf{x}$  is an *eigenvector* and  $\lambda$  is the associated *eigenvalue*

- $\mathbf{A}$  stretches the eigenvector  $\mathbf{x}$
- $\mathbf{A}$  stretches  $\mathbf{x}$  by  $\lambda$
- To find eigenvectors/values: (eigen in  $\mathbb{R}$ )
  - Find  $\lambda$  that solves  $\det(\mathbf{A} - \lambda\mathbf{I}) = 0$
  - Find vectors in **null space** of:

# An Introduction to Eigenvectors, Values, and Diagonalization

## Definition

Suppose  $\mathbf{A}$  is an  $N \times N$  matrix and  $\lambda$  is a scalar.

If

$$\mathbf{Ax} = \lambda\mathbf{x}$$

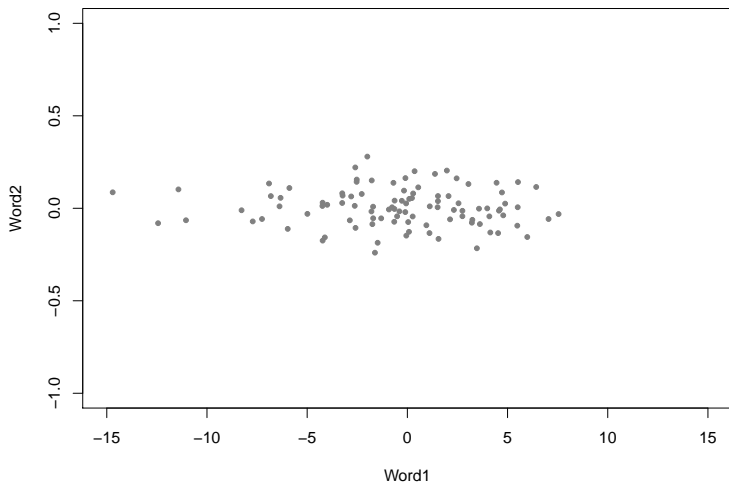
Then  $\mathbf{x}$  is an *eigenvector* and  $\lambda$  is the associated *eigenvalue*

- $\mathbf{A}$  stretches the eigenvector  $\mathbf{x}$
- $\mathbf{A}$  stretches  $\mathbf{x}$  by  $\lambda$
- To find eigenvectors/values: (eigen in  $\mathbb{R}$ )
  - Find  $\lambda$  that solves  $\det(\mathbf{A} - \lambda\mathbf{I}) = 0$
  - Find vectors in **null space** of:

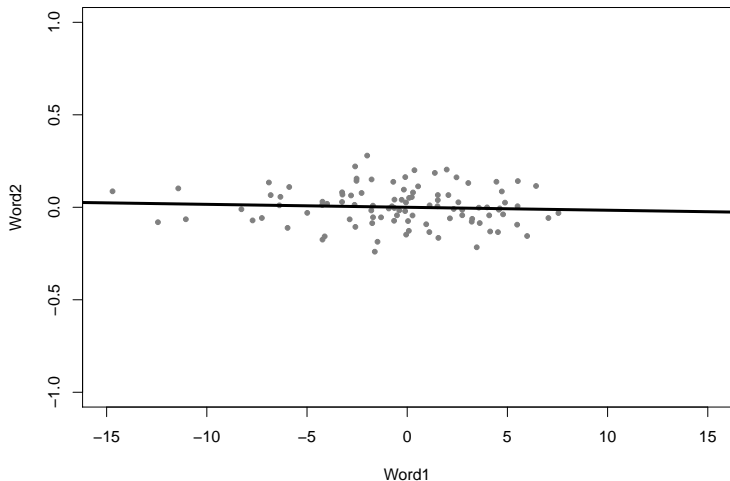
$$(\mathbf{A} - \lambda\mathbf{I}) = 0$$



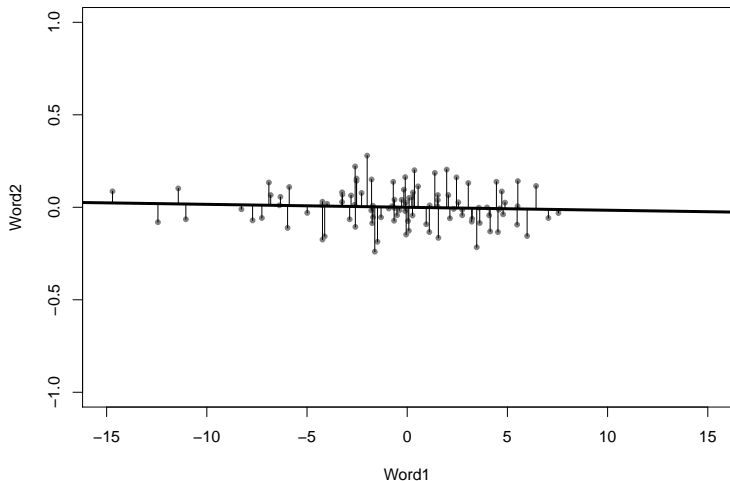
# Finding a Lower Dimensional Space (Manifold Learning)



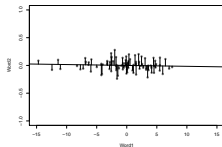
# Finding a Lower Dimensional Space (Manifold Learning)



# Finding a Lower Dimensional Space (Manifold Learning)

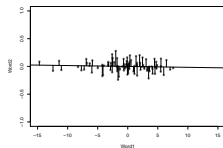


# Finding a Lower Dimensional Space (Manifold Learning)



Original data:

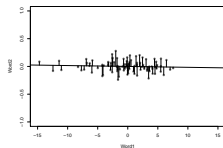
# Finding a Lower Dimensional Space (Manifold Learning)



Original data:

$$\mathbf{x}_i = (x_{i1}, x_{i2})$$

# Finding a Lower Dimensional Space (Manifold Learning)

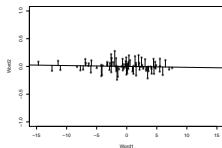


Original data:

$$\mathbf{x}_i = (x_{i1}, x_{i2})$$

Which we approximate with

# Finding a Lower Dimensional Space (Manifold Learning)



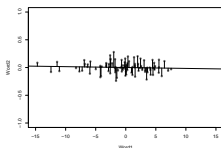
Original data:

$$\mathbf{x}_i = (x_{i1}, x_{i2})$$

Which we approximate with

$$\begin{aligned}\tilde{\mathbf{x}}_i &= z_i \mathbf{w}_1 \\ &= z_i (w_{11}, w_{12})\end{aligned}$$

# Finding a Lower Dimensional Space (Manifold Learning)



Original data  $\mathbf{x}_i \in \mathbb{R}^J$

$$\mathbf{x}_i = (x_{i1}, x_{i2}, \dots, x_{iJ})$$

Which we approximate with  $L (< J)$  weights  $z_{il}$  and vectors  $\mathbf{w}_l \in \mathbb{R}^J$

$$\tilde{\mathbf{x}}_i = z_{i1}\mathbf{w}_1 + z_{i2}\mathbf{w}_2 + \dots + z_{iL}\mathbf{w}_L$$

Define  $\theta = (\underbrace{\mathbf{Z}}_{N \times L}, \underbrace{\mathbf{W}_L}_{L \times J})$



# Principal Component Analysis $\rightsquigarrow$ Objective function

Consider 1-dimensional case ( $L = 1$ ), centered data, and  $\|\mathbf{w}_1\| = 1$ .

# Principal Component Analysis $\rightsquigarrow$ Objective function

Consider 1-dimensional case ( $L = 1$ ), centered data, and  $\|\mathbf{w}_1\| = 1$ .

$$f(\boldsymbol{\theta}, \mathbf{X}) = \frac{1}{N} \sum_{i=1}^N \|\mathbf{x}_i - z_{i1} \mathbf{w}_1\|^2$$

# Principal Component Analysis $\rightsquigarrow$ Objective function

Consider 1-dimensional case ( $L = 1$ ), centered data, and  $\|\mathbf{w}_1\| = 1$ .

$$\begin{aligned} f(\boldsymbol{\theta}, \mathbf{X}) &= \frac{1}{N} \sum_{i=1}^N \|\mathbf{x}_i - z_{i1} \mathbf{w}_1\|^2 \\ &= \frac{1}{N} \sum_{i=1}^N (\mathbf{x}_i - z_{i1} \mathbf{w}_1)' (\mathbf{x}_i - z_{i1} \mathbf{w}_1) \end{aligned}$$

# Principal Component Analysis $\rightsquigarrow$ Objective function

Consider 1-dimensional case ( $L = 1$ ), centered data, and  $\|\mathbf{w}_1\| = 1$ .

$$\begin{aligned} f(\boldsymbol{\theta}, \mathbf{X}) &= \frac{1}{N} \sum_{i=1}^N \|\mathbf{x}_i - z_{i1} \mathbf{w}_1\|^2 \\ &= \frac{1}{N} \sum_{i=1}^N (\mathbf{x}_i - z_{i1} \mathbf{w}_1)' (\mathbf{x}_i - z_{i1} \mathbf{w}_1) \\ &= \frac{1}{N} \sum_{i=1}^N \left( \mathbf{x}_i' \mathbf{x}_i - 2z_{i1} \mathbf{w}_1' \mathbf{x}_i + z_{i1}^2 \right) \end{aligned}$$

# Principal Component Analysis $\rightsquigarrow$ Objective function

Consider 1-dimensional case ( $L = 1$ ), centered data, and  $\|\mathbf{w}_1\| = 1$ .

$$\begin{aligned} f(\boldsymbol{\theta}, \mathbf{X}) &= \frac{1}{N} \sum_{i=1}^N \|\mathbf{x}_i - z_{i1} \mathbf{w}_1\|^2 \\ &= \frac{1}{N} \sum_{i=1}^N (\mathbf{x}_i - z_{i1} \mathbf{w}_1)' (\mathbf{x}_i - z_{i1} \mathbf{w}_1) \\ &= \frac{1}{N} \sum_{i=1}^N \left( \mathbf{x}_i' \mathbf{x}_i - 2z_{i1} \mathbf{w}_1' \mathbf{x}_i + z_{i1}^2 \right) \end{aligned}$$

$$\mathbf{w}_1' \mathbf{w}_1 = 1$$

# Principal Component Analysis $\rightsquigarrow$ Optimization

Optimization:

# Principal Component Analysis $\rightsquigarrow$ Optimization

Optimization:

$$\frac{\partial f(\boldsymbol{\theta}, \mathbf{X})}{\partial z_{i1}} = -\frac{2\mathbf{w}'_1 \mathbf{x}_i + 2z_{i1}}{N}$$

# Principal Component Analysis $\rightsquigarrow$ Optimization

Optimization:

$$\begin{aligned}\frac{\partial f(\boldsymbol{\theta}, \mathbf{X})}{\partial z_{i1}} &= -\frac{2\mathbf{w}'_1 \mathbf{x}_i + 2z_{i1}}{N} \\ 0 &= -\frac{2\mathbf{w}'_1 \mathbf{x}_i + 2z_{i1}^*}{N}\end{aligned}$$



# Principal Component Analysis $\rightsquigarrow$ Optimization

Optimization:

$$\begin{aligned}\frac{\partial f(\boldsymbol{\theta}, \mathbf{X})}{\partial z_{i1}} &= -\frac{2\mathbf{w}'_1 \mathbf{x}_i + 2z_{i1}}{N} \\ 0 &= -\frac{2\mathbf{w}'_1 \mathbf{x}_i + 2z_{i1}^*}{N} \\ z_{i1}^* &= \mathbf{w}'_1 \mathbf{x}_i\end{aligned}$$

# Principal Component Analysis $\rightsquigarrow$ Optimization

Substituting in  $z_{i1}^*$

# Principal Component Analysis $\rightsquigarrow$ Optimization

Substituting in  $z_{i1}^*$

$$= \frac{1}{N} \sum_{i=1}^N (\mathbf{x}_i - z_{i1}^* \mathbf{w}_1)' (\mathbf{x}_i - z_{i1}^* \mathbf{w}_1)$$

# Principal Component Analysis $\rightsquigarrow$ Optimization

Substituting in  $z_{i1}^*$

$$\begin{aligned} &= \frac{1}{N} \sum_{i=1}^N (\mathbf{x}_i - z_{i1}^* \mathbf{w}_1)' (\mathbf{x}_i - z_{i1}^* \mathbf{w}_1) \\ &= \frac{1}{N} \sum_{i=1}^N \left( \underbrace{\mathbf{x}_i' \mathbf{x}_i}_{\text{Constant}} - 2z_{i1}^* \underbrace{\mathbf{w}_1' \mathbf{x}_i}_{z_{i1}^*} + (z_{i1}^*)^2 \underbrace{\mathbf{w}_1' \mathbf{w}_1}_1 \right) \end{aligned}$$

# Principal Component Analysis $\rightsquigarrow$ Optimization

Substituting in  $z_{i1}^*$

$$\begin{aligned} &= \frac{1}{N} \sum_{i=1}^N (\mathbf{x}_i - z_{i1}^* \mathbf{w}_1)' (\mathbf{x}_i - z_{i1}^* \mathbf{w}_1) \\ &= \frac{1}{N} \sum_{i=1}^N \left( \underbrace{\mathbf{x}_i' \mathbf{x}_i}_{\text{Constant}} - 2z_{i1}^* \underbrace{\mathbf{w}_1' \mathbf{x}_i}_{z_{i1}^*} + (z_{i1}^*)^2 \underbrace{\mathbf{w}_1' \mathbf{w}_1}_1 \right) \\ &= -\frac{1}{N} \sum_{i=1}^N (z_{i1}^*)^2 + c \end{aligned}$$

# Principal Component Analysis $\rightsquigarrow$ Optimization

Substituting in  $z_{i1}^*$

$$\begin{aligned} &= \frac{1}{N} \sum_{i=1}^N (\mathbf{x}_i - z_{i1}^* \mathbf{w}_1)' (\mathbf{x}_i - z_{i1}^* \mathbf{w}_1) \\ &= \frac{1}{N} \sum_{i=1}^N \left( \underbrace{\mathbf{x}_i' \mathbf{x}_i}_{\text{Constant}} - 2z_{i1}^* \underbrace{\mathbf{w}_1' \mathbf{x}_i}_{z_{i1}^*} + (z_{i1}^*)^2 \underbrace{\mathbf{w}_1' \mathbf{w}_1}_1 \right) \\ &= -\frac{1}{N} \sum_{i=1}^N (z_{i1}^*)^2 + c \\ &= -\frac{1}{N} \sum_{i=1}^N \mathbf{w}_1' \mathbf{x}_i \mathbf{x}_i' \mathbf{w}_1 \end{aligned}$$

# Principal Component Analysis $\rightsquigarrow$ Optimization

Substituting in  $z_{i1}^*$

$$\begin{aligned} &= \frac{1}{N} \sum_{i=1}^N (\mathbf{x}_i - z_{i1}^* \mathbf{w}_1)' (\mathbf{x}_i - z_{i1}^* \mathbf{w}_1) \\ &= \frac{1}{N} \sum_{i=1}^N \left( \underbrace{\mathbf{x}_i' \mathbf{x}_i}_{\text{Constant}} - 2z_{i1}^* \underbrace{\mathbf{w}_1' \mathbf{x}_i}_{z_{i1}^*} + (z_{i1}^*)^2 \underbrace{\mathbf{w}_1' \mathbf{w}_1}_1 \right) \\ &= -\frac{1}{N} \sum_{i=1}^N (z_{i1}^*)^2 + c \\ &= -\frac{1}{N} \sum_{i=1}^N \mathbf{w}_1' \mathbf{x}_i \mathbf{x}_i' \mathbf{w}_1 \\ &= -\mathbf{w}_1' \boldsymbol{\Sigma} \mathbf{w}_1 \end{aligned}$$

# Principal Component Analysis $\rightsquigarrow$ Optimization

$$= - \sum_{i=1}^N \mathbf{w}_1' \boldsymbol{\Sigma} \mathbf{w}_1$$



# Principal Component Analysis $\rightsquigarrow$ Optimization

$$= - \sum_{i=1}^N \mathbf{w}_1' \boldsymbol{\Sigma} \mathbf{w}_1$$

where  $\boldsymbol{\Sigma}$  is the :

# Principal Component Analysis $\rightsquigarrow$ Optimization

$$= - \sum_{i=1}^N \mathbf{w}_1' \boldsymbol{\Sigma} \mathbf{w}_1$$

where  $\boldsymbol{\Sigma}$  is the :

- Empirical covariance matrix  $\rightsquigarrow \frac{1}{N} \mathbf{X}' \mathbf{X}$

# Principal Component Analysis $\rightsquigarrow$ Optimization

$$= - \sum_{i=1}^N \mathbf{w}_1' \boldsymbol{\Sigma} \mathbf{w}_1$$

where  $\boldsymbol{\Sigma}$  is the :

- Empirical covariance matrix  $\rightsquigarrow \frac{1}{N} \mathbf{X}' \mathbf{X}$
- **Variance** of the projected data. Define

# Principal Component Analysis $\rightsquigarrow$ Optimization

$$= - \sum_{i=1}^N \mathbf{w}_1' \boldsymbol{\Sigma} \mathbf{w}_1$$

where  $\boldsymbol{\Sigma}$  is the :

- Empirical covariance matrix  $\rightsquigarrow \frac{1}{N} \mathbf{X}' \mathbf{X}$
- **Variance** of the projected data. Define

$$\mathbf{z}_1 = (\mathbf{w}_1 \mathbf{x}_1, \mathbf{w}_1 \mathbf{x}_2, \dots, \mathbf{w}_1 \mathbf{x}_N)$$

# Principal Component Analysis $\rightsquigarrow$ Optimization

$$= - \sum_{i=1}^N \mathbf{w}_1' \boldsymbol{\Sigma} \mathbf{w}_1$$

where  $\boldsymbol{\Sigma}$  is the :

- Empirical covariance matrix  $\rightsquigarrow \frac{1}{N} \mathbf{X}' \mathbf{X}$
- **Variance** of the projected data. Define

$$\begin{aligned} \mathbf{z}_1 &= (\mathbf{w}_1 \mathbf{x}_1, \mathbf{w}_1 \mathbf{x}_2, \dots, \mathbf{w}_1 \mathbf{x}_N) \\ \text{var}(\mathbf{z}_1) &= E[\mathbf{z}_1^2] - E[\mathbf{z}_1]^2 \end{aligned}$$

# Principal Component Analysis $\rightsquigarrow$ Optimization

$$= - \sum_{i=1}^N \mathbf{w}'_1 \boldsymbol{\Sigma} \mathbf{w}_1$$

where  $\boldsymbol{\Sigma}$  is the :

- Empirical covariance matrix  $\rightsquigarrow \frac{1}{N} \mathbf{X}' \mathbf{X}$
- **Variance** of the projected data. Define

$$\begin{aligned} \mathbf{z}_1 &= (\mathbf{w}_1 \mathbf{x}_1, \mathbf{w}_1 \mathbf{x}_2, \dots, \mathbf{w}_1 \mathbf{x}_N) \\ \text{var}(\mathbf{z}_1) &= E[\mathbf{z}_1^2] - E[\mathbf{z}_1]^2 \\ &= \frac{1}{N} \sum_{i=1}^N z_{i1}^2 - 0 \end{aligned}$$

# Principal Component Analysis $\rightsquigarrow$ Optimization

$$= - \sum_{i=1}^N \mathbf{w}_1' \boldsymbol{\Sigma} \mathbf{w}_1$$

where  $\boldsymbol{\Sigma}$  is the :

- Empirical covariance matrix  $\rightsquigarrow \frac{1}{N} \mathbf{X}' \mathbf{X}$
- **Variance** of the projected data. Define

$$\begin{aligned} \mathbf{z}_1 &= (\mathbf{w}_1 \mathbf{x}_1, \mathbf{w}_1 \mathbf{x}_2, \dots, \mathbf{w}_1 \mathbf{x}_N) \\ \text{var}(\mathbf{z}_1) &= E[\mathbf{z}_1^2] - E[\mathbf{z}_1]^2 \\ &= \frac{1}{N} \sum_{i=1}^N z_{i1}^2 - 0 \\ &= \frac{1}{N} \sum_{i=1}^N \mathbf{w}_1' \mathbf{x}_i \mathbf{x}_i' \mathbf{w}_1 \end{aligned}$$

# Principal Component Analysis $\rightsquigarrow$ Optimization

$$= - \sum_{i=1}^N \mathbf{w}_1' \boldsymbol{\Sigma} \mathbf{w}_1$$

where  $\boldsymbol{\Sigma}$  is the :

- Empirical covariance matrix  $\rightsquigarrow \frac{1}{N} \mathbf{X}' \mathbf{X}$
- **Variance** of the projected data. Define

$$\begin{aligned} \mathbf{z}_1 &= (\mathbf{w}_1 \mathbf{x}_1, \mathbf{w}_1 \mathbf{x}_2, \dots, \mathbf{w}_1 \mathbf{x}_N) \\ \text{var}(\mathbf{z}_1) &= E[\mathbf{z}_1^2] - E[\mathbf{z}_1]^2 \\ &= \frac{1}{N} \sum_{i=1}^N z_{i1}^2 - 0 \\ &= \frac{1}{N} \sum_{i=1}^N \mathbf{w}_1' \mathbf{x}_i \mathbf{x}_i' \mathbf{w}_1 \end{aligned}$$

Minimize reconstruction error



# Principal Component Analysis $\rightsquigarrow$ Optimization

$$= - \sum_{i=1}^N \mathbf{w}_1' \boldsymbol{\Sigma} \mathbf{w}_1$$

where  $\boldsymbol{\Sigma}$  is the :

- Empirical covariance matrix  $\rightsquigarrow \frac{1}{N} \mathbf{X}' \mathbf{X}$
- **Variance** of the projected data. Define

$$\begin{aligned} \mathbf{z}_1 &= (\mathbf{w}_1 \mathbf{x}_1, \mathbf{w}_1 \mathbf{x}_2, \dots, \mathbf{w}_1 \mathbf{x}_N) \\ \text{var}(\mathbf{z}_1) &= E[\mathbf{z}_1^2] - E[\mathbf{z}_1]^2 \\ &= \frac{1}{N} \sum_{i=1}^N z_{i1}^2 - 0 \\ &= \frac{1}{N} \sum_{i=1}^N \mathbf{w}_1' \mathbf{x}_i \mathbf{x}_i' \mathbf{w}_1 \end{aligned}$$

Minimize reconstruction error  $\rightsquigarrow$  maximize variance of projected data

# Principal Component Analysis $\rightsquigarrow$ Optimization

Maximize variance, subject to constraints

# Principal Component Analysis $\rightsquigarrow$ Optimization

Maximize variance, subject to constraints

$$g(\mathbf{z}^*, \mathbf{w}_1, \mathbf{X}) = \mathbf{w}'_1 \boldsymbol{\Sigma} \mathbf{w}_1 - \lambda_1 (\mathbf{w}'_1 \mathbf{w}_1 - 1)$$

# Principal Component Analysis $\rightsquigarrow$ Optimization

Maximize variance, subject to constraints

$$\begin{aligned}g(\mathbf{z}^*, \mathbf{w}_1, \mathbf{X}) &= \mathbf{w}'_1 \boldsymbol{\Sigma} \mathbf{w}_1 - \lambda_1 (\mathbf{w}'_1 \mathbf{w}_1 - 1) \\ \frac{\partial g(\mathbf{z}^*, \mathbf{w}_1, \mathbf{X})}{\partial \mathbf{w}_1} &= 2\boldsymbol{\Sigma} \mathbf{w}_1 - 2\lambda_1 \mathbf{w}_1\end{aligned}$$

# Principal Component Analysis $\rightsquigarrow$ Optimization

Maximize variance, subject to constraints

$$\begin{aligned}g(\mathbf{z}^*, \mathbf{w}_1, \mathbf{X}) &= \mathbf{w}'_1 \boldsymbol{\Sigma} \mathbf{w}_1 - \lambda_1 (\mathbf{w}'_1 \mathbf{w}_1 - 1) \\ \frac{\partial g(\mathbf{z}^*, \mathbf{w}_1, \mathbf{X})}{\partial \mathbf{w}_1} &= 2\boldsymbol{\Sigma} \mathbf{w}_1 - 2\lambda_1 \mathbf{w}_1 \\ \boldsymbol{\Sigma} \mathbf{w}_1^* &= \lambda_1 \mathbf{w}_1^*\end{aligned}$$

# Principal Component Analysis $\rightsquigarrow$ Optimization

Maximize variance, subject to constraints

$$\begin{aligned}g(\mathbf{z}^*, \mathbf{w}_1, \mathbf{X}) &= \mathbf{w}'_1 \boldsymbol{\Sigma} \mathbf{w}_1 - \lambda_1 (\mathbf{w}'_1 \mathbf{w}_1 - 1) \\ \frac{\partial g(\mathbf{z}^*, \mathbf{w}_1, \mathbf{X})}{\partial \mathbf{w}_1} &= 2\boldsymbol{\Sigma} \mathbf{w}_1 - 2\lambda_1 \mathbf{w}_1 \\ \boldsymbol{\Sigma} \mathbf{w}_1^* &= \lambda_1 \mathbf{w}_1^*\end{aligned}$$

$\mathbf{w}_1^*$  = Eigenvector of  $\boldsymbol{\Sigma}$

# Principal Component Analysis $\rightsquigarrow$ Optimization

Maximize variance, subject to constraints

$$\begin{aligned}g(\mathbf{z}^*, \mathbf{w}_1, \mathbf{X}) &= \mathbf{w}'_1 \boldsymbol{\Sigma} \mathbf{w}_1 - \lambda_1 (\mathbf{w}'_1 \mathbf{w}_1 - 1) \\ \frac{\partial g(\mathbf{z}^*, \mathbf{w}_1, \mathbf{X})}{\partial \mathbf{w}_1} &= 2\boldsymbol{\Sigma} \mathbf{w}_1 - 2\lambda_1 \mathbf{w}_1 \\ \boldsymbol{\Sigma} \mathbf{w}_1^* &= \lambda_1 \mathbf{w}_1^*\end{aligned}$$

$\mathbf{w}_1^*$  = Eigenvector of  $\boldsymbol{\Sigma}$  (!!!!!!!)

# Principal Component Analysis $\rightsquigarrow$ Optimization

Maximize variance, subject to constraints

$$\begin{aligned}g(\mathbf{z}^*, \mathbf{w}_1, \mathbf{X}) &= \mathbf{w}'_1 \boldsymbol{\Sigma} \mathbf{w}_1 - \lambda_1 (\mathbf{w}'_1 \mathbf{w}_1 - 1) \\ \frac{\partial g(\mathbf{z}^*, \mathbf{w}_1, \mathbf{X})}{\partial \mathbf{w}_1} &= 2\boldsymbol{\Sigma} \mathbf{w}_1 - 2\lambda_1 \mathbf{w}_1 \\ \boldsymbol{\Sigma} \mathbf{w}_1^* &= \lambda_1 \mathbf{w}_1^*\end{aligned}$$

$\mathbf{w}_1^*$  = Eigenvector of  $\boldsymbol{\Sigma}$  (!!!!!!!)

We want  $\mathbf{w}_1$  to maximize variance and



# Principal Component Analysis $\rightsquigarrow$ Optimization

Maximize variance, subject to constraints

$$\begin{aligned}g(\mathbf{z}^*, \mathbf{w}_1, \mathbf{X}) &= \mathbf{w}'_1 \boldsymbol{\Sigma} \mathbf{w}_1 - \lambda_1 (\mathbf{w}'_1 \mathbf{w}_1 - 1) \\ \frac{\partial g(\mathbf{z}^*, \mathbf{w}_1, \mathbf{X})}{\partial \mathbf{w}_1} &= 2\boldsymbol{\Sigma} \mathbf{w}_1 - 2\lambda_1 \mathbf{w}_1 \\ \boldsymbol{\Sigma} \mathbf{w}_1^* &= \lambda_1 \mathbf{w}_1^*\end{aligned}$$

$\mathbf{w}_1^*$  = Eigenvector of  $\boldsymbol{\Sigma}$  (!!!!!!!)

We want  $\mathbf{w}_1$  to maximize variance and

$$\mathbf{w}'_1 \boldsymbol{\Sigma} \mathbf{w}_1 = \lambda_1$$

# Principal Component Analysis $\rightsquigarrow$ Optimization

Maximize variance, subject to constraints

$$\begin{aligned}g(\mathbf{z}^*, \mathbf{w}_1, \mathbf{X}) &= \mathbf{w}'_1 \boldsymbol{\Sigma} \mathbf{w}_1 - \lambda_1 (\mathbf{w}'_1 \mathbf{w}_1 - 1) \\ \frac{\partial g(\mathbf{z}^*, \mathbf{w}_1, \mathbf{X})}{\partial \mathbf{w}_1} &= 2\boldsymbol{\Sigma} \mathbf{w}_1 - 2\lambda_1 \mathbf{w}_1 \\ \boldsymbol{\Sigma} \mathbf{w}_1^* &= \lambda_1 \mathbf{w}_1^*\end{aligned}$$

$\mathbf{w}_1^*$  = Eigenvector of  $\boldsymbol{\Sigma}$  (!!!!!!!)

We want  $\mathbf{w}_1$  to maximize variance and

$$\mathbf{w}'_1 \boldsymbol{\Sigma} \mathbf{w}_1 = \lambda_1$$

So  $\mathbf{w}_1$  is eigenvector associated with the largest eigenvalue  $\lambda_1$

# An Introduction to Eigenvectors, Values, and Diagonalization

## Theorem

Suppose  $\mathbf{A}$  is an *invertible*  $N \times N$  matrix. Then  $\mathbf{A}$  has  $N$  distinct eigenvalues and  $N$  linearly independent eigenvectors. Further, we can write  $\mathbf{A}$  as,

$$\mathbf{A} = \mathbf{W}' \begin{pmatrix} \lambda_1 & 0 & \dots & 0 \\ 0 & \lambda_2 & \dots & 0 \\ \vdots & \vdots & \ddots & \vdots \\ 0 & 0 & \dots & \lambda_N \end{pmatrix} \mathbf{W}$$

where  $\mathbf{W} = (\mathbf{w}_1, \mathbf{w}_2, \dots, \mathbf{w}_N)$  is an  $N \times N$  matrix with the  $N$  eigenvectors as column vectors.

# An Introduction to Eigenvectors, Values, and Diagonalization

## Definition

*Suppose  $A$  is a covariance matrix. Then, we can write  $A$  as*

$$\mathbf{A} = \mathbf{W}' \begin{pmatrix} \lambda_1 & 0 & \dots & 0 \\ 0 & \lambda_2 & \dots & 0 \\ \vdots & \vdots & \ddots & \vdots \\ 0 & 0 & \dots & \lambda_N \end{pmatrix} \mathbf{W}$$

*Where  $\lambda_1 > \lambda_2 > \dots > \lambda_N \geq 0$ .*

*We will call  $\mathbf{w}_1$  the first eigenvector,  $\mathbf{w}_2$  the second eigenvector, ...,  $\mathbf{w}_j$  the  $j^{\text{th}}$  eigenvector.*

# Back to Principal Components

## Theorem

*Suppose we want to approximate  $N$  observations  $\mathbf{x}_i \in \mathbb{R}^J$  with  $L < J$  orthogonal-unit length vectors  $\mathbf{w}_l \in \mathbb{R}^J$  with associated scores  $z_{il}$  to minimize reconstruction error:*

# Back to Principal Components

## Theorem

Suppose we want to approximate  $N$  observations  $\mathbf{x}_i \in \mathbb{R}^J$  with  $L < J$  orthogonal-unit length vectors  $\mathbf{w}_l \in \mathbb{R}^J$  with associated scores  $z_{il}$  to minimize reconstruction error:

$$f(\mathbf{X}, \boldsymbol{\theta}) = \frac{1}{N} \sum_{i=1}^N \left\| \mathbf{x}_i - \sum_{l=1}^L z_{il} \mathbf{w}_l \right\|^2$$

# Back to Principal Components

## Theorem

*Suppose we want to approximate  $N$  observations  $\mathbf{x}_i \in \mathbb{R}^J$  with  $L < J$  orthogonal-unit length vectors  $\mathbf{w}_l \in \mathbb{R}^J$  with associated scores  $z_{il}$  to minimize reconstruction error:*

$$f(\mathbf{X}, \boldsymbol{\theta}) = \frac{1}{N} \sum_{i=1}^N \left\| \mathbf{x}_i - \sum_{l=1}^L z_{il} \mathbf{w}_l \right\|^2$$

*The optimal solution sets each  $\mathbf{w}_l$  to be the  $l^{\text{th}}$  eigenvector of the empirical covariance matrix.*

# Back to Principal Components

## Theorem

*Suppose we want to approximate  $N$  observations  $\mathbf{x}_i \in \mathbb{R}^J$  with  $L < J$  orthogonal-unit length vectors  $\mathbf{w}_l \in \mathbb{R}^J$  with associated scores  $z_{il}$  to minimize reconstruction error:*

$$f(\mathbf{X}, \theta) = \frac{1}{N} \sum_{i=1}^N \left\| \mathbf{x}_i - \sum_{l=1}^L z_{il} \mathbf{w}_l \right\|^2$$

*The optimal solution sets each  $\mathbf{w}_l$  to be the  $l^{\text{th}}$  eigenvector of the empirical covariance matrix. Further  $z_{il}^* = \mathbf{w}_l' \mathbf{x}_i$  so that the  $L$  dimensional representation is:*



# Back to Principal Components

## Theorem

Suppose we want to approximate  $N$  observations  $\mathbf{x}_i \in \mathbb{R}^J$  with  $L < J$  orthogonal-unit length vectors  $\mathbf{w}_l \in \mathbb{R}^J$  with associated scores  $z_{il}$  to minimize reconstruction error:

$$f(\mathbf{X}, \theta) = \frac{1}{N} \sum_{i=1}^N \left\| \mathbf{x}_i - \sum_{l=1}^L z_{il} \mathbf{w}_l \right\|^2$$

The optimal solution sets each  $\mathbf{w}_l$  to be the  $l^{\text{th}}$  eigenvector of the empirical covariance matrix. Further  $z_{il}^* = \mathbf{w}_l' \mathbf{x}_i$  so that the  $L$  dimensional representation is:

$$\mathbf{z}_i^L = (\mathbf{w}'_1 \mathbf{x}_i, \mathbf{w}'_2 \mathbf{x}_i, \dots, \mathbf{w}'_L \mathbf{x}_i)$$

# Application of Principal Components in R

Consider press releases from 2005 US Senators

# Application of Principal Components in R

Consider press releases from 2005 US Senators

Define  $\mathbf{x}_i = (x_{i1}, x_{i2}, \dots, x_{iJ})$  as the rate senator  $i$  uses  $J$  words.

# Application of Principal Components in R

Consider press releases from 2005 US Senators

Define  $\mathbf{x}_i = (x_{i1}, x_{i2}, \dots, x_{iJ})$  as the rate senator  $i$  uses  $J$  words.

$$x_{ij} = \frac{\text{No. Times } i \text{ uses word } j}{\text{No. words } i \text{ uses}}$$

# Application of Principal Components in R

Consider press releases from 2005 US Senators

Define  $\mathbf{x}_i = (x_{i1}, x_{i2}, \dots, x_{iJ})$  as the rate senator  $i$  uses  $J$  words.

$$x_{ij} = \frac{\text{No. Times } i \text{ uses word } j}{\text{No. words } i \text{ uses}}$$

dtm:  $100 \times 2796$  matrix containing word rates for senators

# Application of Principal Components in R

Consider press releases from 2005 US Senators

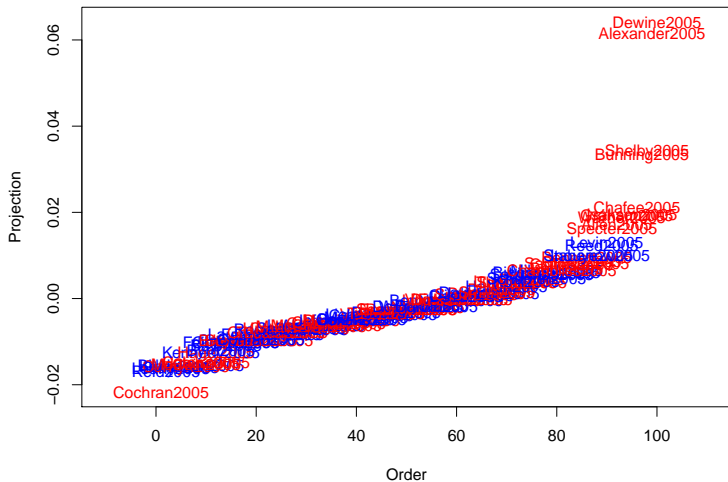
Define  $\mathbf{x}_i = (x_{i1}, x_{i2}, \dots, x_{iJ})$  as the rate senator  $i$  uses  $J$  words.

$$x_{ij} = \frac{\text{No. Times } i \text{ uses word } j}{\text{No. words } i \text{ uses}}$$

dtm:  $100 \times 2796$  matrix containing word rates for senators

`prcomp(dtm)` applies principal components

# Application of Principal Components in R







How do we select the number of dimensions  $L? \rightsquigarrow$  **Model**

We want to minimize reconstruction error

How do we select the number of dimensions  $L$ ?  $\rightsquigarrow$  **Model**

We want to minimize reconstruction error  $\rightsquigarrow$  how well did we do?

How do we select the number of dimensions  $L$ ?  $\rightsquigarrow$  Model

We want to minimize reconstruction error  $\rightsquigarrow$  how well did we do?

$$\text{error}(L) = \frac{1}{N} \sum_{i=1}^N \left\| \mathbf{x}_i - \sum_{l=1}^L z_{il} \mathbf{w}_l \right\|^2$$

How do we select the number of dimensions  $L$ ?  $\rightsquigarrow$  Model

We want to minimize reconstruction error  $\rightsquigarrow$  how well did we do?

$$\text{error}(L) = \frac{1}{N} \sum_{i=1}^N \left\| \mathbf{x}_i - \sum_{l=1}^L z_{il} \mathbf{w}_l \right\|^2$$

Simplifying:

# How do we select the number of dimensions $L$ ? $\rightsquigarrow$ Model

We want to minimize reconstruction error  $\rightsquigarrow$  how well did we do?

$$\text{error}(L) = \frac{1}{N} \sum_{i=1}^N \left\| \mathbf{x}_i - \sum_{l=1}^L z_{il} \mathbf{w}_l \right\|^2$$

Simplifying:

$$\text{error}(L) = \frac{1}{N} \sum_{i=1}^N \left( \mathbf{x}_i - \sum_{l=1}^L z_{il} \mathbf{w}_l \right)' \left( \mathbf{x}_i - \sum_{l=1}^L z_{il} \mathbf{w}_l \right)$$

# How do we select the number of dimensions $L$ ? $\rightsquigarrow$ Model

We want to minimize reconstruction error  $\rightsquigarrow$  how well did we do?

$$\text{error}(L) = \frac{1}{N} \sum_{i=1}^N \left\| \mathbf{x}_i - \sum_{l=1}^L z_{il} \mathbf{w}_l \right\|^2$$

Simplifying:

$$\text{error}(L) = \frac{1}{N} \sum_{i=1}^N \left( \mathbf{x}_i - \sum_{l=1}^L z_{il} \mathbf{w}_l \right)' \left( \mathbf{x}_i - \sum_{l=1}^L z_{il} \mathbf{w}_l \right)$$

Four types of terms:

# How do we select the number of dimensions $L$ ? $\rightsquigarrow$ Model

We want to minimize reconstruction error  $\rightsquigarrow$  how well did we do?

$$\text{error}(L) = \frac{1}{N} \sum_{i=1}^N \left\| \mathbf{x}_i - \sum_{l=1}^L z_{il} \mathbf{w}_l \right\|^2$$

Simplifying:

$$\text{error}(L) = \frac{1}{N} \sum_{i=1}^N \left( \mathbf{x}_i - \sum_{l=1}^L z_{il} \mathbf{w}_l \right)' \left( \mathbf{x}_i - \sum_{l=1}^L z_{il} \mathbf{w}_l \right)$$

Four types of terms:

1)  $\mathbf{x}_i' \mathbf{x}_i$

# How do we select the number of dimensions $L$ ? $\rightsquigarrow$ Model

We want to minimize reconstruction error  $\rightsquigarrow$  how well did we do?

$$\text{error}(L) = \frac{1}{N} \sum_{i=1}^N \left\| \mathbf{x}_i - \sum_{l=1}^L z_{il} \mathbf{w}_l \right\|^2$$

Simplifying:

$$\text{error}(L) = \frac{1}{N} \sum_{i=1}^N \left( \mathbf{x}_i - \sum_{l=1}^L z_{il} \mathbf{w}_l \right)' \left( \mathbf{x}_i - \sum_{l=1}^L z_{il} \mathbf{w}_l \right)$$

Four types of terms:

- 1)  $\mathbf{x}_i' \mathbf{x}_i$
- 2)  $z_{ij} z_{ik} \mathbf{w}_j' \mathbf{w}_k = z_{ij} z_{ik} 0 = 0$



# How do we select the number of dimensions $L$ ? $\rightsquigarrow$ Model

We want to minimize reconstruction error  $\rightsquigarrow$  how well did we do?

$$\text{error}(L) = \frac{1}{N} \sum_{i=1}^N \left\| \mathbf{x}_i - \sum_{l=1}^L z_{il} \mathbf{w}_l \right\|^2$$

Simplifying:

$$\text{error}(L) = \frac{1}{N} \sum_{i=1}^N \left( \mathbf{x}_i - \sum_{l=1}^L z_{il} \mathbf{w}_l \right)' \left( \mathbf{x}_i - \sum_{l=1}^L z_{il} \mathbf{w}_l \right)$$

Four types of terms:

- 1)  $\mathbf{x}_i' \mathbf{x}_i$
- 2)  $z_{ij} z_{ik} \mathbf{w}_j' \mathbf{w}_k = z_{ij} z_{ik} 0 = 0$
- 3)  $z_{ij} z_{ij} \mathbf{w}_j' \mathbf{w}_j = z_{ij}^2$

# How do we select the number of dimensions $L$ ? $\rightsquigarrow$ Model

We want to minimize reconstruction error  $\rightsquigarrow$  how well did we do?

$$\text{error}(L) = \frac{1}{N} \sum_{i=1}^N \left\| \mathbf{x}_i - \sum_{l=1}^L z_{il} \mathbf{w}_l \right\|^2$$

Simplifying:

$$\text{error}(L) = \frac{1}{N} \sum_{i=1}^N \left( \mathbf{x}_i - \sum_{l=1}^L z_{il} \mathbf{w}_l \right)' \left( \mathbf{x}_i - \sum_{l=1}^L z_{il} \mathbf{w}_l \right)$$

Four types of terms:

- 1)  $\mathbf{x}_i' \mathbf{x}_i$
- 2)  $z_{ij} z_{ik} \mathbf{w}_j' \mathbf{w}_k = z_{ij} z_{ik} 0 = 0$
- 3)  $z_{ij} z_{ij} \mathbf{w}_j' \mathbf{w}_j = z_{ij}^2$
- 4)  $\mathbf{x}_i' \sum_{l=1}^L z_{il} \mathbf{w}_l = \sum_{l=1}^L z_{il}^2$

# How do we select the number of dimensions $L$ ? $\rightsquigarrow$ Model

We want to minimize reconstruction error  $\rightsquigarrow$  how well did we do?

$$\text{error}(L) = \frac{1}{N} \sum_{i=1}^N \left\| \mathbf{x}_i - \sum_{l=1}^L z_{il} \mathbf{w}_l \right\|^2$$

Simplifying:

$$\begin{aligned} \text{error}(L) &= \frac{1}{N} \sum_{i=1}^N (\mathbf{x}_i - \sum_{l=1}^L z_{il} \mathbf{w}_l)' (\mathbf{x}_i - \sum_{l=1}^L z_{il} \mathbf{w}_l) \\ &= \frac{1}{N} \sum_{i=1}^N \left( \mathbf{x}_i' \mathbf{x}_i - \sum_{l=1}^L z_{il}^2 \right) \end{aligned}$$

Four types of terms:

- 1)  $\mathbf{x}_i' \mathbf{x}_i$
- 2)  $z_{ij} z_{ik} \mathbf{w}_j' \mathbf{w}_k = z_{ij} z_{ik} 0 = 0$
- 3)  $z_{ij} z_{ij} \mathbf{w}_j' \mathbf{w}_j = z_{ij}^2$
- 4)  $\mathbf{x}_i' \sum_{l=1}^L z_{il} \mathbf{w}_l = \sum_{l=1}^L z_{il}^2$

How do we select the number of dimensions  $L$ ?  $\rightsquigarrow$  **Model**

$$\text{error}(L) = \frac{1}{N} \sum_{i=1}^N \left( \mathbf{x}'_i \mathbf{x}_i - \sum_{l=1}^L z_{il}^2 \right)$$

How do we select the number of dimensions  $L$ ?  $\rightsquigarrow$  **Model**

$$\begin{aligned}\text{error}(L) &= \frac{1}{N} \sum_{i=1}^N \left( \mathbf{x}'_i \mathbf{x}_i - \sum_{l=1}^L z_{il}^2 \right) \\ &= \frac{1}{N} \sum_{i=1}^N \left( \mathbf{x}'_i \mathbf{x}_i - \sum_{l=1}^L \mathbf{w}_l \mathbf{x}_i \mathbf{x}'_i \mathbf{w}_l \right)\end{aligned}$$

How do we select the number of dimensions  $L$ ?  $\rightsquigarrow$  **Model**

$$\begin{aligned}\text{error}(L) &= \frac{1}{N} \sum_{i=1}^N \left( \mathbf{x}'_i \mathbf{x}_i - \sum_{l=1}^L z_{il}^2 \right) \\ &= \frac{1}{N} \sum_{i=1}^N \left( \mathbf{x}'_i \mathbf{x}_i - \sum_{l=1}^L \mathbf{w}_l \mathbf{x}_i \mathbf{x}'_i \mathbf{w}_l \right) \\ &= \frac{1}{N} \sum_{i=1}^N (\mathbf{x}'_i \mathbf{x}_i) - \frac{1}{N} \sum_{l=1}^L \sum_{i=1}^N \mathbf{w}'_l \mathbf{x}_i \mathbf{x}'_i \mathbf{w}_l\end{aligned}$$

How do we select the number of dimensions  $L$ ?  $\rightsquigarrow$  Model

$$\begin{aligned}\text{error}(L) &= \frac{1}{N} \sum_{i=1}^N \left( \mathbf{x}'_i \mathbf{x}_i - \sum_{l=1}^L z_{il}^2 \right) \\ &= \frac{1}{N} \sum_{i=1}^N \left( \mathbf{x}'_i \mathbf{x}_i - \sum_{l=1}^L \mathbf{w}_l \mathbf{x}_i \mathbf{x}'_i \mathbf{w}_l \right) \\ &= \frac{1}{N} \sum_{i=1}^N (\mathbf{x}'_i \mathbf{x}_i) - \frac{1}{N} \sum_{l=1}^L \sum_{i=1}^N \mathbf{w}'_l \mathbf{x}_i \mathbf{x}'_i \mathbf{w}_l \\ &= \frac{1}{N} \sum_{i=1}^N (\mathbf{x}'_i \mathbf{x}_i) - \sum_{l=1}^L \mathbf{w}'_l \boldsymbol{\Sigma} \mathbf{w}_l\end{aligned}$$

How do we select the number of dimensions  $L$ ?  $\rightsquigarrow$  **Model**

$$\begin{aligned}\text{error}(L) &= \frac{1}{N} \sum_{i=1}^N \left( \mathbf{x}'_i \mathbf{x}_i - \sum_{l=1}^L z_{il}^2 \right) \\ &= \frac{1}{N} \sum_{i=1}^N \left( \mathbf{x}'_i \mathbf{x}_i - \sum_{l=1}^L \mathbf{w}_l \mathbf{x}_i \mathbf{x}'_i \mathbf{w}'_l \right) \\ &= \frac{1}{N} \sum_{i=1}^N (\mathbf{x}'_i \mathbf{x}_i) - \frac{1}{N} \sum_{l=1}^L \sum_{i=1}^N \mathbf{w}'_l \mathbf{x}_i \mathbf{x}'_i \mathbf{w}_l \\ &= \frac{1}{N} \sum_{i=1}^N (\mathbf{x}'_i \mathbf{x}_i) - \sum_{l=1}^L \mathbf{w}'_l \boldsymbol{\Sigma} \mathbf{w}_l \\ &= \frac{1}{N} \sum_{i=1}^N (\mathbf{x}'_i \mathbf{x}_i) - \sum_{l=1}^L \lambda_l \mathbf{w}'_l \mathbf{w}_l\end{aligned}$$



How do we select the number of dimensions  $L$ ?  $\rightsquigarrow$  Model

$$\begin{aligned}\text{error}(L) &= \frac{1}{N} \sum_{i=1}^N \left( \mathbf{x}'_i \mathbf{x}_i - \sum_{l=1}^L z_{il}^2 \right) \\ &= \frac{1}{N} \sum_{i=1}^N \left( \mathbf{x}'_i \mathbf{x}_i - \sum_{l=1}^L \mathbf{w}_l \mathbf{x}_i \mathbf{x}'_i \mathbf{w}'_l \right) \\ &= \frac{1}{N} \sum_{i=1}^N (\mathbf{x}'_i \mathbf{x}_i) - \frac{1}{N} \sum_{l=1}^L \sum_{i=1}^N \mathbf{w}'_l \mathbf{x}_i \mathbf{x}'_i \mathbf{w}_l \\ &= \frac{1}{N} \sum_{i=1}^N (\mathbf{x}'_i \mathbf{x}_i) - \sum_{l=1}^L \mathbf{w}'_l \boldsymbol{\Sigma} \mathbf{w}_l \\ &= \frac{1}{N} \sum_{i=1}^N (\mathbf{x}'_i \mathbf{x}_i) - \sum_{l=1}^L \lambda_l \mathbf{w}'_l \mathbf{w}_l \\ &= \frac{1}{N} \sum_{i=1}^N (\mathbf{x}'_i \mathbf{x}_i) - \sum_{l=1}^L \lambda_l\end{aligned}$$

How do we select the number of dimensions  $L$ ?  $\rightsquigarrow$  Model

If  $L = J$

How do we select the number of dimensions  $L$ ?  $\rightsquigarrow$  Model

If  $L = J$

$$\text{error}(J) = \frac{1}{N} \sum_{i=1}^N (\mathbf{x}'_i \mathbf{x}_i) - \sum_{l=1}^J \lambda_l = 0$$

How do we select the number of dimensions  $L$ ?  $\rightsquigarrow$  Model

If  $L = J$

$$\text{error}(J) = \frac{1}{N} \sum_{i=1}^N (\mathbf{x}'_i \mathbf{x}_i) - \sum_{l=1}^J \lambda_l = 0$$

So for  $L < J$ ,

How do we select the number of dimensions  $L$ ?  $\rightsquigarrow$  Model

If  $L = J$

$$\text{error}(J) = \frac{1}{N} \sum_{i=1}^N (\mathbf{x}'_i \mathbf{x}_i) - \sum_{l=1}^J \lambda_l = 0$$

So for  $L < J$ ,

$$0 = \frac{1}{N} \sum_{i=1}^N (\mathbf{x}'_i \mathbf{x}_i) - \left( \sum_{l=1}^L \lambda_l + \sum_{j=L+1}^J \lambda_j \right)$$

How do we select the number of dimensions  $L$ ?  $\rightsquigarrow$  Model

If  $L = J$

$$\text{error}(J) = \frac{1}{N} \sum_{i=1}^N (\mathbf{x}'_i \mathbf{x}_i) - \sum_{l=1}^J \lambda_l = 0$$

So for  $L < J$ ,

$$0 = \frac{1}{N} \sum_{i=1}^N (\mathbf{x}'_i \mathbf{x}_i) - \left( \sum_{l=1}^L \lambda_l + \sum_{j=L+1}^J \lambda_j \right)$$

$$\sum_{j=L+1}^J \lambda_j = \frac{1}{N} \sum_{i=1}^N (\mathbf{x}'_i \mathbf{x}_i) - \sum_{l=1}^L \lambda_l$$

How do we select the number of dimensions  $L$ ?  $\rightsquigarrow$  Model

If  $L = J$

$$\text{error}(J) = \frac{1}{N} \sum_{i=1}^N (\mathbf{x}'_i \mathbf{x}_i) - \sum_{l=1}^J \lambda_l = 0$$

So for  $L < J$ ,

$$0 = \frac{1}{N} \sum_{i=1}^N (\mathbf{x}'_i \mathbf{x}_i) - \left( \sum_{l=1}^L \lambda_l + \sum_{j=L+1}^J \lambda_j \right)$$

$$\sum_{j=L+1}^J \lambda_j = \frac{1}{N} \sum_{i=1}^N (\mathbf{x}'_i \mathbf{x}_i) - \sum_{l=1}^L \lambda_l$$

$$\sum_{j=L+1}^J \lambda_j = \text{error}(L)$$

How do we select the number of dimensions  $L$ ?  $\rightsquigarrow$  Model

$$\sum_{j=L+1}^J \lambda_j = \text{error}(L)$$



How do we select the number of dimensions  $L$ ?  $\rightsquigarrow$  Model

$$\sum_{j=L+1}^J \lambda_j = \text{error}(L)$$

- Error = Sum of “remaining” eigenvalues

How do we select the number of dimensions  $L$ ?  $\rightsquigarrow$  Model

$$\sum_{j=L+1}^J \lambda_j = \text{error}(L)$$

- Error = Sum of “remaining” eigenvalues
- Total variance explained = (sum of included eigenvalues)/(sum of all eigenvalues)

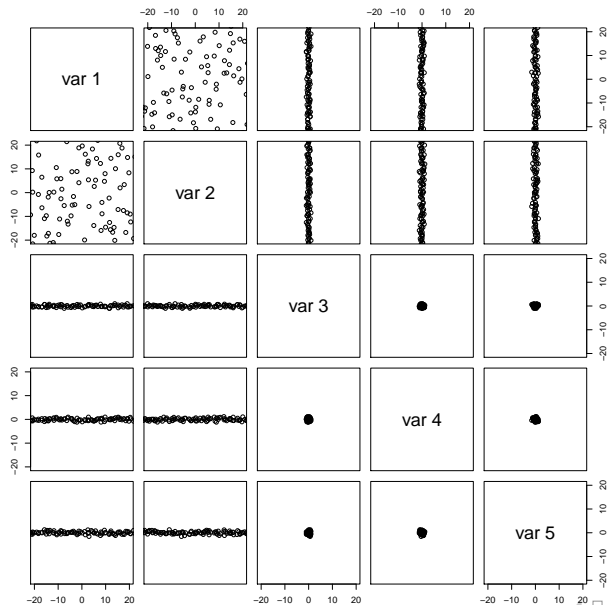
How do we select the number of dimensions  $L$ ?  $\rightsquigarrow$  Model

$$\sum_{j=L+1}^J \lambda_j = \text{error}(L)$$

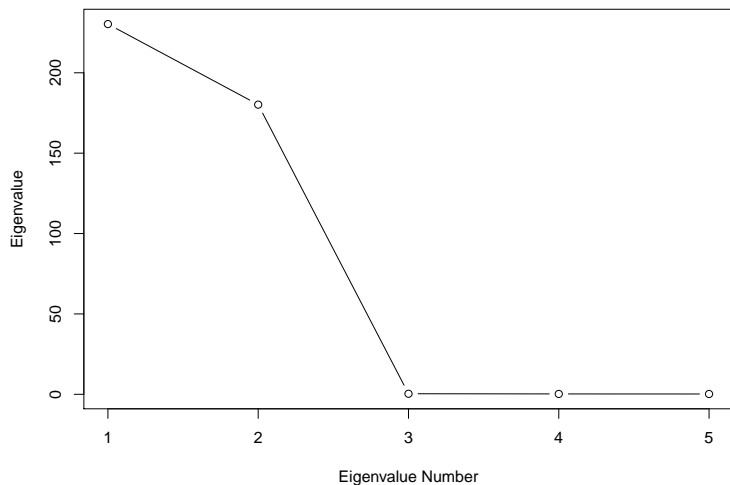
- Error = Sum of “remaining” eigenvalues
- Total variance explained = (sum of included eigenvalues)/(sum of all eigenvalues)

Recommendation  $\rightsquigarrow$  look for Elbow

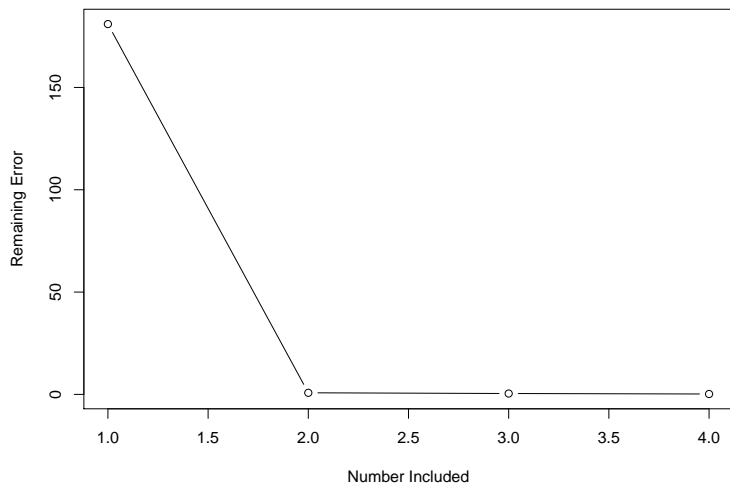
# How do we select the number of dimensions $L$ ? $\rightsquigarrow$ Model



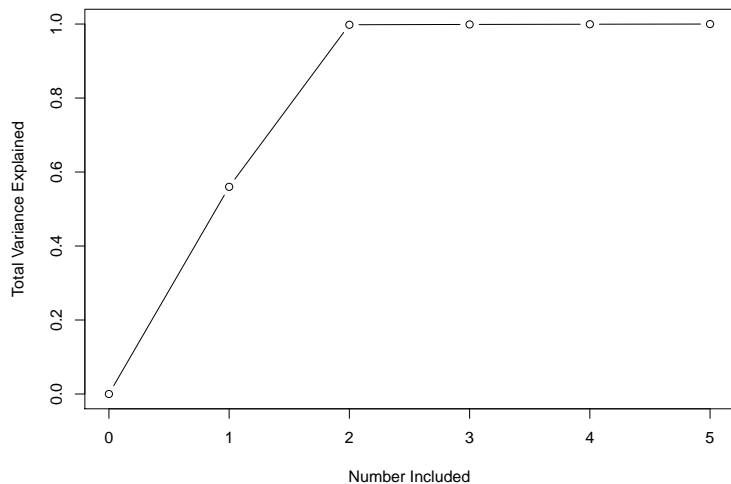
How do we select the number of dimensions  $L$ ?  $\rightsquigarrow$  Model



How do we select the number of dimensions  $L$ ?  $\rightsquigarrow$  Model



How do we select the number of dimensions  $L$ ?  $\rightsquigarrow$  Model



# Non-model based evaluations: What's the point?

What is the true underlying dimensionality of  $\mathbf{X}$ ?



# Non-model based evaluations: What's the point?

What is the true underlying dimensionality of  $\mathbf{X}$ ? **J**

# Non-model based evaluations: What's the point?

What is the true underlying dimensionality of  $\mathbf{X}$ ?  $J$ (!!!!!!)

# Non-model based evaluations: What's the point?

What is the true underlying dimensionality of  $\mathbf{X}$ ?  $J$ (!!!!!!)

- Attempts to assess dimensionality require a **model**  $\rightsquigarrow$  some way to tradeoff accuracy of reconstruction with simplicity

# Non-model based evaluations: What's the point?

What is the true underlying dimensionality of  $\mathbf{X}$ ?  $J$ (!!!!!!)

- Attempts to assess dimensionality require a **model**  $\rightsquigarrow$  some way to tradeoff accuracy of reconstruction with simplicity
- **Any** answer (no matter how creatively obtained) supposes **you have the right function to measure tradeoff**

# Non-model based evaluations: What's the point?

What is the true underlying dimensionality of  $\mathbf{X}$ ?  $J$ (!!!!!!)

- Attempts to assess dimensionality require a **model**  $\rightsquigarrow$  some way to tradeoff accuracy of reconstruction with simplicity
- **Any** answer (no matter how creatively obtained) supposes **you have the right function to measure tradeoff**
- The “right” number of dimensions depends on the **task** you have in mind

# Non-model based evaluations: What's the point?

What is the true underlying dimensionality of  $\mathbf{X}$ ?  $J$ (!!!!!!)

- Attempts to assess dimensionality require a **model**  $\rightsquigarrow$  some way to tradeoff accuracy of reconstruction with simplicity
- **Any** answer (no matter how creatively obtained) supposes **you have the right function to measure tradeoff**
- The “right” number of dimensions depends on the **task** you have in mind

Mathematical model  $\rightsquigarrow$  insufficient to make modeling decision

# Kernel Principal Component Analysis

Recall from Thursday that we define a **Kernel** ( $N \times N$ ) matrix as:

$$\mathbf{K} = \begin{pmatrix} k(\mathbf{x}_1, \mathbf{x}_1) & k(\mathbf{x}_1, \mathbf{x}_2) & \dots & k(\mathbf{x}_1, \mathbf{x}_N) \\ k(\mathbf{x}_2, \mathbf{x}_1) & k(\mathbf{x}_2, \mathbf{x}_2) & \dots & k(\mathbf{x}_2, \mathbf{x}_N) \\ \vdots & \vdots & \ddots & \vdots \\ k(\mathbf{x}_N, \mathbf{x}_1) & k(\mathbf{x}_N, \mathbf{x}_2) & \dots & k(\mathbf{x}_N, \mathbf{x}_N) \end{pmatrix}$$

# Kernel Principal Component Analysis

Recall from Thursday that we define a **Kernel** ( $N \times N$ ) matrix as:

$$\mathbf{K} = \begin{pmatrix} k(\mathbf{x}_1, \mathbf{x}_1) & k(\mathbf{x}_1, \mathbf{x}_2) & \dots & k(\mathbf{x}_1, \mathbf{x}_N) \\ k(\mathbf{x}_2, \mathbf{x}_1) & k(\mathbf{x}_2, \mathbf{x}_2) & \dots & k(\mathbf{x}_2, \mathbf{x}_N) \\ \vdots & \vdots & \ddots & \vdots \\ k(\mathbf{x}_N, \mathbf{x}_1) & k(\mathbf{x}_N, \mathbf{x}_2) & \dots & k(\mathbf{x}_N, \mathbf{x}_N) \end{pmatrix}$$

Where we suppose this matrix emerges from applying  $\phi : \mathbb{R}^J \rightarrow \mathbb{R}^M$  to the data and then taking the inner product:



# Kernel Principal Component Analysis

Recall from Thursday that we define a **Kernel** ( $N \times N$ ) matrix as:

$$\mathbf{K} = \begin{pmatrix} k(\mathbf{x}_1, \mathbf{x}_1) & k(\mathbf{x}_1, \mathbf{x}_2) & \dots & k(\mathbf{x}_1, \mathbf{x}_N) \\ k(\mathbf{x}_2, \mathbf{x}_1) & k(\mathbf{x}_2, \mathbf{x}_2) & \dots & k(\mathbf{x}_2, \mathbf{x}_N) \\ \vdots & \vdots & \ddots & \vdots \\ k(\mathbf{x}_N, \mathbf{x}_1) & k(\mathbf{x}_N, \mathbf{x}_2) & \dots & k(\mathbf{x}_N, \mathbf{x}_N) \end{pmatrix}$$

Where we suppose this matrix emerges from applying  $\phi : \mathbb{R}^J \rightarrow \mathbb{R}^M$  to the data and then taking the inner product:

$$\mathbf{K} = \mathbf{\Phi}\mathbf{\Phi}' \text{ (The inner product matrix)}$$

# Kernel Principal Component Analysis

Recall from Thursday that we define a **Kernel** ( $N \times N$ ) matrix as:

$$\mathbf{K} = \begin{pmatrix} k(\mathbf{x}_1, \mathbf{x}_1) & k(\mathbf{x}_1, \mathbf{x}_2) & \dots & k(\mathbf{x}_1, \mathbf{x}_N) \\ k(\mathbf{x}_2, \mathbf{x}_1) & k(\mathbf{x}_2, \mathbf{x}_2) & \dots & k(\mathbf{x}_2, \mathbf{x}_N) \\ \vdots & \vdots & \ddots & \vdots \\ k(\mathbf{x}_N, \mathbf{x}_1) & k(\mathbf{x}_N, \mathbf{x}_2) & \dots & k(\mathbf{x}_N, \mathbf{x}_N) \end{pmatrix}$$

Where we suppose this matrix emerges from applying  $\phi : \mathbb{R}^J \rightarrow \mathbb{R}^M$  to the data and then taking the inner product:

$$\begin{aligned} \mathbf{K} &= \mathbf{\Phi}\mathbf{\Phi}' \text{ (The inner product matrix)} \\ &= \begin{pmatrix} \phi(\mathbf{x}_1)' \phi(\mathbf{x}_1) & \phi(\mathbf{x}_1)' \phi(\mathbf{x}_2) & \dots & \phi(\mathbf{x}_1)' \phi(\mathbf{x}_N) \\ \phi(\mathbf{x}_2)' \phi(\mathbf{x}_1) & \phi(\mathbf{x}_2)' \phi(\mathbf{x}_2) & \dots & \phi(\mathbf{x}_2)' \phi(\mathbf{x}_N) \\ \vdots & \vdots & \ddots & \vdots \\ \phi(\mathbf{x}_N)' \phi(\mathbf{x}_1) & \phi(\mathbf{x}_N)' \phi(\mathbf{x}_2) & \dots & \phi(\mathbf{x}_N)' \phi(\mathbf{x}_N) \end{pmatrix} \end{aligned}$$

# Kernel Principal Component Analysis

Recall from Thursday that we define a **Kernel** ( $N \times N$ ) matrix as:

$$\mathbf{K} = \begin{pmatrix} k(\mathbf{x}_1, \mathbf{x}_1) & k(\mathbf{x}_1, \mathbf{x}_2) & \dots & k(\mathbf{x}_1, \mathbf{x}_N) \\ k(\mathbf{x}_2, \mathbf{x}_1) & k(\mathbf{x}_2, \mathbf{x}_2) & \dots & k(\mathbf{x}_2, \mathbf{x}_N) \\ \vdots & \vdots & \ddots & \vdots \\ k(\mathbf{x}_N, \mathbf{x}_1) & k(\mathbf{x}_N, \mathbf{x}_2) & \dots & k(\mathbf{x}_N, \mathbf{x}_N) \end{pmatrix}$$

Where we suppose this matrix emerges from applying  $\phi : \mathbb{R}^J \rightarrow \mathbb{R}^M$  to the data and then taking the inner product:

$$\begin{aligned} \mathbf{K} &= \mathbf{\Phi}\mathbf{\Phi}' \text{ (The inner product matrix)} \\ &= \begin{pmatrix} \phi(\mathbf{x}_1)' \phi(\mathbf{x}_1) & \phi(\mathbf{x}_1)' \phi(\mathbf{x}_2) & \dots & \phi(\mathbf{x}_1)' \phi(\mathbf{x}_N) \\ \phi(\mathbf{x}_2)' \phi(\mathbf{x}_1) & \phi(\mathbf{x}_2)' \phi(\mathbf{x}_2) & \dots & \phi(\mathbf{x}_2)' \phi(\mathbf{x}_N) \\ \vdots & \vdots & \ddots & \vdots \\ \phi(\mathbf{x}_N)' \phi(\mathbf{x}_1) & \phi(\mathbf{x}_N)' \phi(\mathbf{x}_2) & \dots & \phi(\mathbf{x}_N)' \phi(\mathbf{x}_N) \end{pmatrix} \end{aligned}$$

Compute PCA of  $\mathbf{\Phi}$  from  $\mathbf{\Phi}\mathbf{\Phi}'$

# Kernel PCA

PCA of  $\mathbf{X}$

# Kernel PCA

PCA of  $\mathbf{X}$  Eigenvectors of  $\mathbf{X}'\mathbf{X}$  ( $\frac{1}{N}$  doesn't affect eigenvectors)

# Kernel PCA

PCA of  $\mathbf{X}$  Eigenvectors of  $\mathbf{X}'\mathbf{X}$  ( $\frac{1}{N}$  doesn't affect eigenvectors)

Suppose  $\mathbf{u}_1$  is an eigenvector for  $\mathbf{X}\mathbf{X}'$ , with value  $\lambda_1$ .

# Kernel PCA

PCA of  $\mathbf{X}$  Eigenvectors of  $\mathbf{X}'\mathbf{X}$  ( $\frac{1}{N}$  doesn't affect eigenvectors)

Suppose  $\mathbf{u}_1$  is an eigenvector for  $\mathbf{X}\mathbf{X}'$ , with value  $\lambda_1$ . Then

# Kernel PCA

PCA of  $\mathbf{X}$  Eigenvectors of  $\mathbf{X}'\mathbf{X}$  ( $\frac{1}{N}$  doesn't affect eigenvectors)

Suppose  $\mathbf{u}_1$  is an eigenvector for  $\mathbf{X}\mathbf{X}'$ , with value  $\lambda_1$ . Then

$$(\mathbf{X}\mathbf{X}')\mathbf{u}_1 = \lambda_1\mathbf{u}_1$$



# Kernel PCA

PCA of  $\mathbf{X}$  Eigenvectors of  $\mathbf{X}'\mathbf{X}$  ( $\frac{1}{N}$  doesn't affect eigenvectors)

Suppose  $\mathbf{u}_1$  is an eigenvector for  $\mathbf{X}\mathbf{X}'$ , with value  $\lambda_1$ . Then

$$\begin{aligned}(\mathbf{X}\mathbf{X}')\mathbf{u}_1 &= \lambda_1\mathbf{u}_1 \\(\mathbf{X}'\mathbf{X})(\mathbf{X}'\mathbf{u}_1) &= \lambda_1(\mathbf{X}'\mathbf{u}_1)\end{aligned}$$

# Kernel PCA

PCA of  $\mathbf{X}$  Eigenvectors of  $\mathbf{X}'\mathbf{X}$  ( $\frac{1}{N}$  doesn't affect eigenvectors)

Suppose  $\mathbf{u}_1$  is an eigenvector for  $\mathbf{X}\mathbf{X}'$ , with value  $\lambda_1$ . Then

$$\begin{aligned}(\mathbf{X}\mathbf{X}')\mathbf{u}_1 &= \lambda_1\mathbf{u}_1 \\(\mathbf{X}'\mathbf{X})(\mathbf{X}'\mathbf{u}_1) &= \lambda_1(\mathbf{X}'\mathbf{u}_1) \\ &= \lambda_1\mathbf{v}_1\end{aligned}$$

# Kernel PCA

PCA of  $\mathbf{X}$  Eigenvectors of  $\mathbf{X}'\mathbf{X}$  ( $\frac{1}{N}$  doesn't affect eigenvectors)

Suppose  $\mathbf{u}_1$  is an eigenvector for  $\mathbf{X}\mathbf{X}'$ , with value  $\lambda_1$ . Then

$$\begin{aligned}(\mathbf{X}\mathbf{X}')\mathbf{u}_1 &= \lambda_1\mathbf{u}_1 \\(\mathbf{X}'\mathbf{X})(\mathbf{X}'\mathbf{u}_1) &= \lambda_1(\mathbf{X}'\mathbf{u}_1) \\ &= \lambda_1\mathbf{v}_1\end{aligned}$$

But  $\mathbf{v}_1$  needs unit length, and

# Kernel PCA

PCA of  $\mathbf{X}$  Eigenvectors of  $\mathbf{X}'\mathbf{X}$  ( $\frac{1}{N}$  doesn't affect eigenvectors)

Suppose  $\mathbf{u}_1$  is an eigenvector for  $\mathbf{X}\mathbf{X}'$ , with value  $\lambda_1$ . Then

$$\begin{aligned}(\mathbf{X}\mathbf{X}')\mathbf{u}_1 &= \lambda_1\mathbf{u}_1 \\(\mathbf{X}'\mathbf{X})(\mathbf{X}'\mathbf{u}_1) &= \lambda_1(\mathbf{X}'\mathbf{u}_1) \\ &= \lambda_1\mathbf{v}_1\end{aligned}$$

But  $\mathbf{v}_1$  needs unit length, and

$$\|\mathbf{v}_1\|^2 = \mathbf{v}_1'\mathbf{v}_1$$

# Kernel PCA

PCA of  $\mathbf{X}$  Eigenvectors of  $\mathbf{X}'\mathbf{X}$  ( $\frac{1}{N}$  doesn't affect eigenvectors)

Suppose  $\mathbf{u}_1$  is an eigenvector for  $\mathbf{X}\mathbf{X}'$ , with value  $\lambda_1$ . Then

$$\begin{aligned}(\mathbf{X}\mathbf{X}')\mathbf{u}_1 &= \lambda_1\mathbf{u}_1 \\(\mathbf{X}'\mathbf{X})(\mathbf{X}'\mathbf{u}_1) &= \lambda_1(\mathbf{X}'\mathbf{u}_1) \\ &= \lambda_1\mathbf{v}_1\end{aligned}$$

But  $\mathbf{v}_1$  needs unit length, and

$$\begin{aligned}\|\mathbf{v}_1\|^2 &= \mathbf{v}_1'\mathbf{v}_1 \\ &= \mathbf{u}_1'\mathbf{X}\mathbf{X}'\mathbf{u}_1\end{aligned}$$

# Kernel PCA

PCA of  $\mathbf{X}$  Eigenvectors of  $\mathbf{X}'\mathbf{X}$  ( $\frac{1}{N}$  doesn't affect eigenvectors)

Suppose  $\mathbf{u}_1$  is an eigenvector for  $\mathbf{X}\mathbf{X}'$ , with value  $\lambda_1$ . Then

$$\begin{aligned}(\mathbf{X}\mathbf{X}')\mathbf{u}_1 &= \lambda_1\mathbf{u}_1 \\(\mathbf{X}'\mathbf{X})(\mathbf{X}'\mathbf{u}_1) &= \lambda_1(\mathbf{X}'\mathbf{u}_1) \\ &= \lambda_1\mathbf{v}_1\end{aligned}$$

But  $\mathbf{v}_1$  needs unit length, and

$$\begin{aligned}\|\mathbf{v}_1\|^2 &= \mathbf{v}_1'\mathbf{v}_1 \\ &= \mathbf{u}_1'\mathbf{X}\mathbf{X}'\mathbf{u}_1 \\ &= \lambda_1\mathbf{u}_1'\mathbf{u}_1 = \lambda_1\end{aligned}$$

# Kernel PCA

PCA of  $\mathbf{X}$  Eigenvectors of  $\mathbf{X}'\mathbf{X}$  ( $\frac{1}{N}$  doesn't affect eigenvectors)

Suppose  $\mathbf{u}_1$  is an eigenvector for  $\mathbf{X}\mathbf{X}'$ , with value  $\lambda_1$ . Then

$$\begin{aligned}(\mathbf{X}\mathbf{X}')\mathbf{u}_1 &= \lambda_1\mathbf{u}_1 \\(\mathbf{X}'\mathbf{X})(\mathbf{X}'\mathbf{u}_1) &= \lambda_1(\mathbf{X}'\mathbf{u}_1) \\ &= \lambda_1\mathbf{v}_1\end{aligned}$$

But  $\mathbf{v}_1$  needs unit length, and

$$\begin{aligned}\|\mathbf{v}_1\|^2 &= \mathbf{v}_1'\mathbf{v}_1 \\ &= \mathbf{u}_1'\mathbf{X}\mathbf{X}'\mathbf{u}_1 \\ &= \lambda_1\mathbf{u}_1'\mathbf{u}_1 = \lambda_1\end{aligned}$$

So first eigenvector of  $\mathbf{X}'\mathbf{X}$  is

# Kernel PCA

PCA of  $\mathbf{X}$  Eigenvectors of  $\mathbf{X}'\mathbf{X}$  ( $\frac{1}{N}$  doesn't affect eigenvectors)

Suppose  $\mathbf{u}_1$  is an eigenvector for  $\mathbf{X}\mathbf{X}'$ , with value  $\lambda_1$ . Then

$$\begin{aligned}(\mathbf{X}\mathbf{X}')\mathbf{u}_1 &= \lambda_1\mathbf{u}_1 \\(\mathbf{X}'\mathbf{X})(\mathbf{X}'\mathbf{u}_1) &= \lambda_1(\mathbf{X}'\mathbf{u}_1) \\ &= \lambda_1\mathbf{v}_1\end{aligned}$$

But  $\mathbf{v}_1$  needs unit length, and

$$\begin{aligned}\|\mathbf{v}_1\|^2 &= \mathbf{v}_1'\mathbf{v}_1 \\ &= \mathbf{u}_1'\mathbf{X}\mathbf{X}'\mathbf{u}_1 \\ &= \lambda_1\mathbf{u}_1'\mathbf{u}_1 = \lambda_1\end{aligned}$$

So first eigenvector of  $\mathbf{X}'\mathbf{X}$  is

$$\mathbf{w}_1 = \frac{1}{\sqrt{\lambda_1}}\mathbf{X}'\mathbf{u}_1$$



# Kernel PCA

$\mathbf{K} = \Phi\Phi'$  (assume  $\Phi$  is mean-centered, for now)

# Kernel PCA

$\mathbf{K} = \mathbf{\Phi}\mathbf{\Phi}'$  (assume  $\mathbf{\Phi}$  is mean-centered, for now)

We can obtain  $\mathbf{u}_1$  and  $\lambda_1$  from  $\mathbf{K}$ .

# Kernel PCA

$\mathbf{K} = \Phi\Phi'$  (assume  $\Phi$  is mean-centered, for now)

We can obtain  $\mathbf{u}_1$  and  $\lambda_1$  from  $\mathbf{K}$ . We know that

# Kernel PCA

$\mathbf{K} = \Phi\Phi'$  (assume  $\Phi$  is mean-centered, for now)

We can obtain  $\mathbf{u}_1$  and  $\lambda_1$  from  $\mathbf{K}$ . We know that

$$\mathbf{w}_1 = \frac{1}{\sqrt{\lambda_1}} \underbrace{\Phi'}_{\text{Unknown}} \mathbf{u}_1$$

# Kernel PCA

$\mathbf{K} = \Phi\Phi'$  (assume  $\Phi$  is mean-centered, for now)

We can obtain  $\mathbf{u}_1$  and  $\lambda_1$  from  $\mathbf{K}$ . We know that

$$\mathbf{w}_1 = \frac{1}{\sqrt{\lambda_1}} \underbrace{\Phi'}_{\text{Unknown}} \mathbf{u}_1$$

But suppose we want to project a point  $\phi(\mathbf{x}_i)$ , then

# Kernel PCA

$\mathbf{K} = \mathbf{\Phi}\mathbf{\Phi}'$  (assume  $\mathbf{\Phi}$  is mean-centered, for now)

We can obtain  $\mathbf{u}_1$  and  $\lambda_1$  from  $\mathbf{K}$ . We know that

$$\mathbf{w}_1 = \frac{1}{\sqrt{\lambda_1}} \underbrace{\mathbf{\Phi}'}_{\text{Unknown}} \mathbf{u}_1$$

But suppose we want to project a point  $\phi(\mathbf{x}_i)$ , then

$$\phi(\mathbf{x}_i)' \mathbf{w}_1 = \frac{1}{\sqrt{\lambda_1}} \phi(\mathbf{x}_i)' \mathbf{\Phi}' \mathbf{u}_1$$

# Kernel PCA

$\mathbf{K} = \mathbf{\Phi}\mathbf{\Phi}'$  (assume  $\mathbf{\Phi}$  is mean-centered, for now)

We can obtain  $\mathbf{u}_1$  and  $\lambda_1$  from  $\mathbf{K}$ . We know that

$$\mathbf{w}_1 = \frac{1}{\sqrt{\lambda_1}} \underbrace{\mathbf{\Phi}'}_{\text{Unknown}} \mathbf{u}_1$$

But suppose we want to project a point  $\phi(\mathbf{x}_i)$ , then

$$\phi(\mathbf{x}_i)' \mathbf{w}_1 = \frac{1}{\sqrt{\lambda_1}} \phi(\mathbf{x}_i)' \mathbf{\Phi}' \mathbf{u}_1$$

$$\phi(\mathbf{x}_i)' \mathbf{\Phi}' = \left[ \phi(\mathbf{x}_i)' \phi(\mathbf{x}_1), \phi(\mathbf{x}_i)' \phi(\mathbf{x}_2), \dots, \phi(\mathbf{x}_i)' \phi(\mathbf{x}_N) \right]$$

# Kernel PCA

$\mathbf{K} = \mathbf{\Phi}\mathbf{\Phi}'$  (assume  $\mathbf{\Phi}$  is mean-centered, for now)

We can obtain  $\mathbf{u}_1$  and  $\lambda_1$  from  $\mathbf{K}$ . We know that

$$\mathbf{w}_1 = \frac{1}{\sqrt{\lambda_1}} \underbrace{\mathbf{\Phi}'}_{\text{Unknown}} \mathbf{u}_1$$

But suppose we want to project a point  $\phi(\mathbf{x}_i)$ , then

$$\phi(\mathbf{x}_i)' \mathbf{w}_1 = \frac{1}{\sqrt{\lambda_1}} \phi(\mathbf{x}_i)' \mathbf{\Phi}' \mathbf{u}_1$$

$$\begin{aligned} \phi(\mathbf{x}_i)' \mathbf{\Phi}' &= \left[ \phi(\mathbf{x}_i)' \phi(\mathbf{x}_1), \phi(\mathbf{x}_i)' \phi(\mathbf{x}_2), \dots, \phi(\mathbf{x}_i)' \phi(\mathbf{x}_N) \right] \\ &= \left[ k(\mathbf{x}_i, \mathbf{x}_1), k(\mathbf{x}_i, \mathbf{x}_2), \dots, k(\mathbf{x}_i, \mathbf{x}_N) \right] \end{aligned}$$



# Kernel PCA

$\mathbf{K} = \mathbf{\Phi}\mathbf{\Phi}'$  (assume  $\mathbf{\Phi}$  is mean-centered, for now)

We can obtain  $\mathbf{u}_1$  and  $\lambda_1$  from  $\mathbf{K}$ . We know that

$$\mathbf{w}_1 = \frac{1}{\sqrt{\lambda_1}} \underbrace{\mathbf{\Phi}'}_{\text{Unknown}} \mathbf{u}_1$$

But suppose we want to project a point  $\phi(\mathbf{x}_i)$ , then

$$\phi(\mathbf{x}_i)' \mathbf{w}_1 = \frac{1}{\sqrt{\lambda_1}} \phi(\mathbf{x}_i)' \mathbf{\Phi}' \mathbf{u}_1$$

$$\begin{aligned} \phi(\mathbf{x}_i)' \mathbf{\Phi}' &= \left[ \phi(\mathbf{x}_i)' \phi(\mathbf{x}_1), \phi(\mathbf{x}_i)' \phi(\mathbf{x}_2), \dots, \phi(\mathbf{x}_i)' \phi(\mathbf{x}_N) \right] \\ &= \left[ k(\mathbf{x}_i, \mathbf{x}_1), k(\mathbf{x}_i, \mathbf{x}_2), \dots, k(\mathbf{x}_i, \mathbf{x}_N) \right] = \mathbf{k}(\mathbf{x}_i, *) \end{aligned}$$

# Kernel PCA

$\mathbf{K} = \mathbf{\Phi}\mathbf{\Phi}'$  (assume  $\mathbf{\Phi}$  is mean-centered, for now)

We can obtain  $\mathbf{u}_1$  and  $\lambda_1$  from  $\mathbf{K}$ . We know that

$$\mathbf{w}_1 = \frac{1}{\sqrt{\lambda_1}} \underbrace{\mathbf{\Phi}'}_{\text{Unknown}} \mathbf{u}_1$$

But suppose we want to project a point  $\phi(\mathbf{x}_i)$ , then

$$\phi(\mathbf{x}_i)' \mathbf{w}_1 = \frac{1}{\sqrt{\lambda_1}} \phi(\mathbf{x}_i)' \mathbf{\Phi}' \mathbf{u}_1$$

$$\begin{aligned} \phi(\mathbf{x}_i)' \mathbf{\Phi}' &= [\phi(\mathbf{x}_i)' \phi(\mathbf{x}_1), \phi(\mathbf{x}_i)' \phi(\mathbf{x}_2), \dots, \phi(\mathbf{x}_i)' \phi(\mathbf{x}_N)] \\ &= [k(\mathbf{x}_i, \mathbf{x}_1), k(\mathbf{x}_i, \mathbf{x}_2), \dots, k(\mathbf{x}_i, \mathbf{x}_N)] = \mathbf{k}(\mathbf{x}_i, *) \end{aligned}$$

Then, we can obtain projection for observation  $i$  using Kernel with

# Kernel PCA

$\mathbf{K} = \mathbf{\Phi}\mathbf{\Phi}'$  (assume  $\mathbf{\Phi}$  is mean-centered, for now)

We can obtain  $\mathbf{u}_1$  and  $\lambda_1$  from  $\mathbf{K}$ . We know that

$$\mathbf{w}_1 = \frac{1}{\sqrt{\lambda_1}} \underbrace{\mathbf{\Phi}'}_{\text{Unknown}} \mathbf{u}_1$$

But suppose we want to project a point  $\phi(\mathbf{x}_i)$ , then

$$\phi(\mathbf{x}_i)' \mathbf{w}_1 = \frac{1}{\sqrt{\lambda_1}} \phi(\mathbf{x}_i)' \mathbf{\Phi}' \mathbf{u}_1$$

$$\begin{aligned} \phi(\mathbf{x}_i)' \mathbf{\Phi}' &= \left[ \phi(\mathbf{x}_i)' \phi(\mathbf{x}_1), \phi(\mathbf{x}_i)' \phi(\mathbf{x}_2), \dots, \phi(\mathbf{x}_i)' \phi(\mathbf{x}_N) \right] \\ &= \left[ k(\mathbf{x}_i, \mathbf{x}_1), k(\mathbf{x}_i, \mathbf{x}_2), \dots, k(\mathbf{x}_i, \mathbf{x}_N) \right] = \mathbf{k}(\mathbf{x}_i, *) \end{aligned}$$

Then, we can obtain projection for observation  $i$  using Kernel with

$$\phi(\mathbf{x}_i)' \mathbf{w}_1 = \frac{1}{\sqrt{\lambda_1}} \mathbf{k}(\mathbf{x}_i, *) \mathbf{u}_1$$

# Kernel PCA

Center  $\mathbf{K}$ ?

Use centering matrix  $\mathbf{H}$

$$\mathbf{H} = \mathbf{I}_N - \frac{(\mathbf{1}_N \mathbf{1}'_N)}{N}$$
$$\mathbf{K}_{\text{center}} = \mathbf{H} \mathbf{K} \mathbf{H}$$

# Spirling and Indian Treaties

**Spirling (2013)**: model Treaties between US and Native Americans

# Spirling and Indian Treaties

**Spirling (2013)**: model Treaties between US and Native Americans  
Why?

# Spirling and Indian Treaties

**Spirling (2013)**: model Treaties between US and Native Americans  
Why?

- American political development

# Spirling and Indian Treaties

**Spirling (2013)**: model Treaties between US and Native Americans  
Why?

- American political development
- IR Theories of Treaties and Treaty Violations



# Spirling and Indian Treaties

**Spirling (2013)**: model Treaties between US and Native Americans  
Why?

- American political development
- IR Theories of Treaties and Treaty Violations
- Comparative studies of indigenous/colonialist interaction

# Spirling and Indian Treaties

**Spirling (2013)**: model Treaties between US and Native Americans  
Why?

- American political development
- IR Theories of Treaties and Treaty Violations
- Comparative studies of indigenous/colonialist interaction
- **Political Science question**: how did Native Americans lose land so quickly?

# Spiraling and Indian Treaties

How do we preserve word order and semantic language?

# Spiraling and Indian Treaties

How do we preserve word order and semantic language?  
After stemming, stopping, bag of wording:

# Spiraling and Indian Treaties

How do we preserve word order and semantic language?

After stemming, stopping, bag of wording:

- Peace Between Us

# Spiraling and Indian Treaties

How do we preserve word order and semantic language?

After stemming, stopping, bag of wording:

- Peace Between Us
- No Peace Between Us

# Spiraling and Indian Treaties

How do we preserve word order and semantic language?

After stemming, stopping, bag of wording:

- Peace Between Us
- No Peace Between Us

are identical.

# Spiraling and Indian Treaties

How do we preserve word order and semantic language?

After stemming, stopping, bag of wording:

- Peace Between Us
- No Peace Between Us

are identical.

Spiraling uses complicated representation of texts to preserve word order



# Spiraling and Indian Treaties

How do we preserve word order and semantic language?

After stemming, stopping, bag of wording:

- Peace Between Us
- No Peace Between Us

are identical.

Spiraling uses complicated representation of texts to preserve word order ~>  
broad application

# Spiraling and Indian Treaties

How do we preserve word order and semantic language?

After stemming, stopping, bag of wording:

- Peace Between Us
- No Peace Between Us

are identical.

Spiraling uses complicated representation of texts to preserve word order↔  
broad application

Peace Between Us

# Spiraling and Indian Treaties

How do we preserve word order and semantic language?

After stemming, stopping, bag of wording:

- Peace Between Us
- No Peace Between Us

are identical.

Spiraling uses complicated representation of texts to preserve word order↔  
broad application

Peace Between Us

# Spiraling and Indian Treaties

How do we preserve word order and semantic language?

After stemming, stopping, bag of wording:

- Peace Between Us
- No Peace Between Us

are identical.

Spiraling uses complicated representation of texts to preserve word order↔  
broad application

Peace Between Us

# Spiraling and Indian Treaties

How do we preserve word order and semantic language?

After stemming, stopping, bag of wording:

- Peace Between Us
- No Peace Between Us

are identical.

Spiraling uses complicated representation of texts to preserve word order↔  
broad application

Peace Between Us

# Spiraling and Indian Treaties

How do we preserve word order and semantic language?

After stemming, stopping, bag of wording:

- Peace Between Us
- No Peace Between Us

are identical.

Spiraling uses complicated representation of texts to preserve word order↔  
broad application

Peace **B**etween Us

# Spiraling and Indian Treaties

How do we preserve word order and semantic language?

After stemming, stopping, bag of wording:

- Peace Between Us
- No Peace Between Us

are identical.

Spiraling uses complicated representation of texts to preserve word order↔  
broad application

Peace **B**etween Us

# Spiraling and Indian Treaties

How do we preserve word order and semantic language?

After stemming, stopping, bag of wording:

- Peace Between Us
- No Peace Between Us

are identical.

Spiraling uses complicated representation of texts to preserve word order↔  
broad application

Peace **B**etween Us



# Spiraling and Indian Treaties

How do we preserve word order and semantic language?

After stemming, stopping, bag of wording:

- Peace Between Us
- No Peace Between Us

are identical.

Spiraling uses complicated representation of texts to preserve word order↔  
broad application

Peace **B**etween Us

# Spiraling and Indian Treaties

How do we preserve word order and semantic language?

After stemming, stopping, bag of wording:

- Peace Between Us
- No Peace Between Us

are identical.

Spiraling uses complicated representation of texts to preserve word order↔  
broad application

Peace **Be**tween Us

# Spiraling and Indian Treaties

How do we preserve word order and semantic language?

After stemming, stopping, bag of wording:

- Peace Between Us
- No Peace Between Us

are identical.

Spiraling uses complicated representation of texts to preserve word order↔  
broad application

Peace Bet**ween** Us

# Spiraling and Indian Treaties

How do we preserve word order and semantic language?

After stemming, stopping, bag of wording:

- Peace Between Us
- No Peace Between Us

are identical.

Spiraling uses complicated representation of texts to preserve word order↔  
broad application

Peace Between **een** Us

# Spiraling and Indian Treaties

How do we preserve word order and semantic language?

After stemming, stopping, bag of wording:

- Peace Between Us
- No Peace Between Us

are identical.

Spiraling uses complicated representation of texts to preserve word order↔  
broad application

Peace Between **en** **U**s

# Spiraling and Indian Treaties

How do we preserve word order and semantic language?

After stemming, stopping, bag of wording:

- Peace Between Us
- No Peace Between Us

are identical.

Spiraling uses complicated representation of texts to preserve word order↔  
broad application

Peace Between **U**s

# Spiraling and Indian Treaties

Consider documents  $x_i$  and  $x_j$ , where we have preserved order, punctuation, and all else.

# Spiraling and Indian Treaties

Consider documents  $\mathbf{x}_i$  and  $\mathbf{x}_j$ , where we have preserved order, punctuation, and all else.

We say  $\mathbf{x}_i \in \mathcal{X}$



# Spiraling and Indian Treaties

Consider documents  $\mathbf{x}_i$  and  $\mathbf{x}_j$ , where we have preserved order, punctuation, and all else.

We say  $\mathbf{x}_i \in \mathcal{X}$

Spiraling examines 5-character strings,  $s \in \mathcal{A}$

# Spiraling and Indian Treaties

Consider documents  $\mathbf{x}_i$  and  $\mathbf{x}_j$ , where we have preserved order, punctuation, and all else.

We say  $\mathbf{x}_i \in \mathcal{X}$

Spiraling examines 5-character strings,  $s \in \mathcal{A}$

Define:

# Spiraling and Indian Treaties

Consider documents  $\mathbf{x}_i$  and  $\mathbf{x}_j$ , where we have preserved order, punctuation, and all else.

We say  $\mathbf{x}_i \in \mathcal{X}$

Spiraling examines 5-character strings,  $s \in \mathcal{A}$

Define:

$\phi_s : \mathcal{X} \rightarrow \mathbb{R}$  as a function that counts the number of times string  $s$  occurs in document  $\mathbf{x}$ .

# Spiraling and Indian Treaties

Consider documents  $\mathbf{x}_i$  and  $\mathbf{x}_j$ , where we have preserved order, punctuation, and all else.

We say  $\mathbf{x}_i \in \mathcal{X}$

Spiraling examines 5-character strings,  $s \in \mathcal{A}$

Define:

$\phi_s : \mathcal{X} \rightarrow \mathbb{R}$  as a function that counts the number of times string  $s$  occurs in document  $\mathbf{x}$ .

Define **string kernel** to be,

# Spiraling and Indian Treaties

Consider documents  $\mathbf{x}_i$  and  $\mathbf{x}_j$ , where we have preserved order, punctuation, and all else.

We say  $\mathbf{x}_i \in \mathcal{X}$

Spiraling examines 5-character strings,  $s \in \mathcal{A}$

Define:

$\phi_s : \mathcal{X} \rightarrow \mathbb{R}$  as a function that counts the number of times string  $s$  occurs in document  $\mathbf{x}$ .

Define **string kernel** to be,

$$k(\mathbf{x}_i, \mathbf{x}_j) = \sum_{s \in \mathcal{A}} w_s \phi_s(\mathbf{x}_i) \phi_s(\mathbf{x}_j)$$

# Spiraling and Indian Treaties

Consider documents  $\mathbf{x}_i$  and  $\mathbf{x}_j$ , where we have preserved order, punctuation, and all else.

We say  $\mathbf{x}_i \in \mathcal{X}$

Spiraling examines 5-character strings,  $s \in \mathcal{A}$

Define:

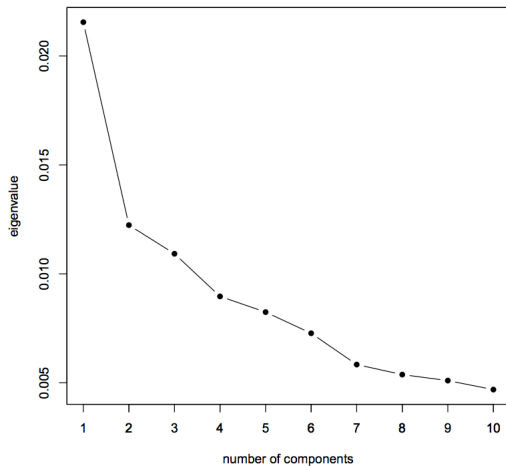
$\phi_s : \mathcal{X} \rightarrow \mathbb{R}$  as a function that counts the number of times string  $s$  occurs in document  $\mathbf{x}$ .

Define **string kernel** to be,

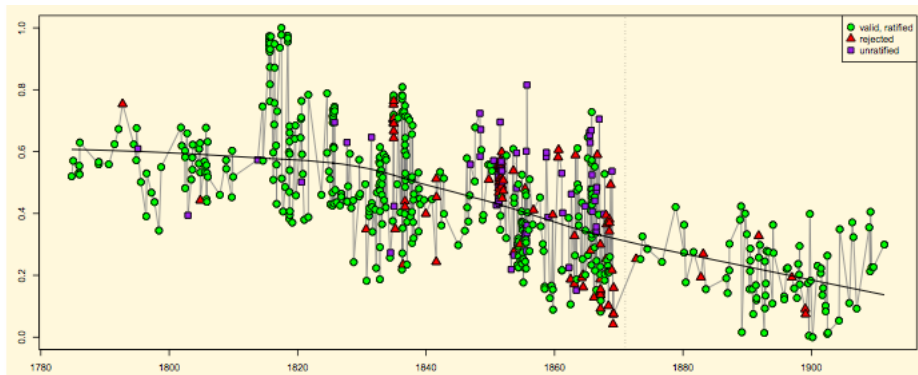
$$k(\mathbf{x}_i, \mathbf{x}_j) = \sum_{s \in \mathcal{A}} w_s \phi_s(\mathbf{x}_i) \phi_s(\mathbf{x}_j)$$

$\phi(\mathbf{x}_i) \approx \binom{32}{5}$  element long count vector

# Spirling and Indian Treaties



# Spirling and Indian Treaties





# Classic Multidimensional Scaling $\rightsquigarrow$ Objective Function

Suppose we have an  $N \times N$  matrix  $\mathbf{D}$   $\rightsquigarrow$

# Classic Multidimensional Scaling $\rightsquigarrow$ Objective Function

Suppose we have an  $N \times N$  matrix  $\mathbf{D}$   $\rightsquigarrow$  matrix of squared Euclidean distances between documents  $\mathbf{x}_i, \mathbf{x}_j \in \mathbb{R}^J$ .

# Classic Multidimensional Scaling $\rightsquigarrow$ Objective Function

Suppose we have an  $N \times N$  matrix  $\mathbf{D}$   $\rightsquigarrow$  matrix of squared Euclidean distances between documents  $\mathbf{x}_i, \mathbf{x}_j \in \mathbb{R}^J$ .

Typical distance between  $\mathbf{x}_i$  and  $\mathbf{x}_j$

# Classic Multidimensional Scaling $\rightsquigarrow$ Objective Function

Suppose we have an  $N \times N$  matrix  $\mathbf{D} \rightsquigarrow$  matrix of squared Euclidean distances between documents  $\mathbf{x}_i, \mathbf{x}_j \in \mathbb{R}^J$ .

Typical distance between  $\mathbf{x}_i$  and  $\mathbf{x}_j$

$$d_{ij}^2 = \sum_{k=1}^J (x_{ik} - x_{jk})^2$$

# Classic Multidimensional Scaling $\rightsquigarrow$ Objective Function

Suppose we have an  $N \times N$  matrix  $\mathbf{D} \rightsquigarrow$  matrix of squared Euclidean distances between documents  $\mathbf{x}_i, \mathbf{x}_j \in \mathbb{R}^J$ .

Typical distance between  $\mathbf{x}_i$  and  $\mathbf{x}_j$

$$d_{ij}^2 = \sum_{k=1}^J (x_{ik} - x_{jk})^2$$

For all  $N$  documents we look for points  $\mathbf{z}_i \in \mathbb{R}^L$  (with  $L < J$ ) to minimize

# Classic Multidimensional Scaling $\rightsquigarrow$ Objective Function

Suppose we have an  $N \times N$  matrix  $\mathbf{D} \rightsquigarrow$  matrix of squared Euclidean distances between documents  $\mathbf{x}_i, \mathbf{x}_j \in \mathbb{R}^J$ .

Typical distance between  $\mathbf{x}_i$  and  $\mathbf{x}_j$

$$d_{ij}^2 = \sum_{k=1}^J (x_{ik} - x_{jk})^2$$

For all  $N$  documents we look for points  $\mathbf{z}_i \in \mathbb{R}^L$  (with  $L < J$ ) to minimize

$$f(\mathbf{X}, \boldsymbol{\theta}) = f(\mathbf{X}, \mathbf{Z})$$

# Classic Multidimensional Scaling $\rightsquigarrow$ Objective Function

Suppose we have an  $N \times N$  matrix  $\mathbf{D} \rightsquigarrow$  matrix of squared Euclidean distances between documents  $\mathbf{x}_i, \mathbf{x}_j \in \mathbb{R}^J$ .

Typical distance between  $\mathbf{x}_i$  and  $\mathbf{x}_j$

$$d_{ij}^2 = \sum_{k=1}^J (x_{ik} - x_{jk})^2$$

For all  $N$  documents we look for points  $\mathbf{z}_i \in \mathbb{R}^L$  (with  $L < J$ ) to minimize

$$\begin{aligned} f(\mathbf{X}, \boldsymbol{\theta}) &= f(\mathbf{X}, \mathbf{Z}) \\ &= \sum_{i \neq j} (d_{ij} - \|\mathbf{z}_i - \mathbf{z}_j\|)^2 \end{aligned}$$

# Classic Multidimensional Scaling $\rightsquigarrow$ Objective Function

Suppose we have an  $N \times N$  matrix  $\mathbf{D} \rightsquigarrow$  matrix of squared Euclidean distances between documents  $\mathbf{x}_i, \mathbf{x}_j \in \mathbb{R}^J$ .

Typical distance between  $\mathbf{x}_i$  and  $\mathbf{x}_j$

$$d_{ij}^2 = \sum_{k=1}^J (x_{ik} - x_{jk})^2$$

For all  $N$  documents we look for points  $\mathbf{z}_i \in \mathbb{R}^L$  (with  $L < J$ ) to minimize

$$\begin{aligned} f(\mathbf{X}, \boldsymbol{\theta}) &= f(\mathbf{X}, \mathbf{Z}) \\ &= \sum_{i \neq j} (d_{ij} - \|\mathbf{z}_i - \mathbf{z}_j\|)^2 \end{aligned}$$

Recover best low-dimensional representation  $\rightsquigarrow$  Eigenvector decomposition of  $\mathbf{D}$



# Classic Multidimensional Scaling $\rightsquigarrow$ Optimization

- 1) Start with  $D$
- 2) “Center”  $D$  to obtain  $XX'$
- 3) Find largest eigenvectors of  $XX'$

# Classic Multidimensional Scaling $\rightsquigarrow$ Optimization

Center the distance matrix, obtain inner product

# Classic Multidimensional Scaling $\rightsquigarrow$ Optimization

Center the distance matrix, obtain inner product

$$d_{ij}^2 = \sum_{k=1}^J (x_{ik} - x_{jk})^2$$

# Classic Multidimensional Scaling $\rightsquigarrow$ Optimization

Center the distance matrix, obtain inner product

$$\begin{aligned}d_{ij}^2 &= \sum_{k=1}^J (x_{ik} - x_{jk})^2 \\ &= \sum_{k=1}^J (x_{ik}^2 - 2x_{ik}x_{jk} + x_{jk}^2)\end{aligned}$$

# Classic Multidimensional Scaling $\rightsquigarrow$ Optimization

Center the distance matrix, obtain inner product

$$\begin{aligned}d_{ij}^2 &= \sum_{k=1}^J (x_{ik} - x_{jk})^2 \\ &= \sum_{k=1}^J (x_{ik}^2 - 2x_{ik}x_{jk} + x_{jk}^2) \\ &= \underbrace{\mathbf{x}'_i \mathbf{x}_i}_{\text{rows}} - 2\mathbf{x}_i \mathbf{x}_j + \underbrace{\mathbf{x}'_j \mathbf{x}_j}_{\text{columns}}\end{aligned}$$

# Classic Multidimensional Scaling $\rightsquigarrow$ Optimization

Center the distance matrix, obtain inner product

$$\begin{aligned}d_{ij}^2 &= \sum_{k=1}^J (x_{ik} - x_{jk})^2 \\ &= \sum_{k=1}^J (x_{ik}^2 - 2x_{ik}x_{jk} + x_{jk}^2) \\ &= \underbrace{\mathbf{x}'_i \mathbf{x}_i}_{\text{rows}} - 2\mathbf{x}_i \mathbf{x}_j + \underbrace{\mathbf{x}'_j \mathbf{x}_j}_{\text{columns}}\end{aligned}$$

- Subtract off row means

# Classic Multidimensional Scaling $\rightsquigarrow$ Optimization

Center the distance matrix, obtain inner product

$$\begin{aligned}d_{ij}^2 &= \sum_{k=1}^J (x_{ik} - x_{jk})^2 \\ &= \sum_{k=1}^J (x_{ik}^2 - 2x_{ik}x_{jk} + x_{jk}^2) \\ &= \underbrace{\mathbf{x}'_i \mathbf{x}_i}_{\text{rows}} - 2\mathbf{x}_i \mathbf{x}_j + \underbrace{\mathbf{x}'_j \mathbf{x}_j}_{\text{columns}}\end{aligned}$$

- Subtract off row means
- Subtract off column means

# Classic Multidimensional Scaling $\rightsquigarrow$ Optimization

Center the distance matrix, obtain inner product

$$\begin{aligned}d_{ij}^2 &= \sum_{k=1}^J (x_{ik} - x_{jk})^2 \\ &= \sum_{k=1}^J (x_{ik}^2 - 2x_{ik}x_{jk} + x_{jk}^2) \\ &= \underbrace{\mathbf{x}'_i \mathbf{x}_i}_{\text{rows}} - 2\mathbf{x}_i \mathbf{x}_j + \underbrace{\mathbf{x}'_j \mathbf{x}_j}_{\text{columns}}\end{aligned}$$

- Subtract off row means
- Subtract off column means
- Divide by -2



# Classic Multidimensional Scaling $\rightsquigarrow$ Optimization

Center the distance matrix, obtain inner product

$$\begin{aligned}d_{ij}^2 &= \sum_{k=1}^J (x_{ik} - x_{jk})^2 \\ &= \sum_{k=1}^J (x_{ik}^2 - 2x_{ik}x_{jk} + x_{jk}^2) \\ &= \underbrace{\mathbf{x}'_i \mathbf{x}_i}_{\text{rows}} - 2\mathbf{x}_i \mathbf{x}_j + \underbrace{\mathbf{x}'_j \mathbf{x}_j}_{\text{columns}}\end{aligned}$$

- Subtract off row means
- Subtract off column means
- Divide by -2

$$\mathbf{X}\mathbf{X}' = \frac{-\mathbf{H}\mathbf{D}\mathbf{H}}{2}$$

# Classic Multidimensional Scaling $\rightsquigarrow$ Optimization

We can write  $\mathbf{X}\mathbf{X}'$  as

$$\begin{aligned}\mathbf{X}\mathbf{X}' &= \mathbf{W}' \begin{pmatrix} \lambda_1 & 0 & \dots & 0 \\ 0 & \lambda_2 & \dots & 0 \\ \vdots & \vdots & \ddots & \vdots \\ 0 & 0 & \dots & \lambda_N \end{pmatrix} \mathbf{W} \\ &= \underbrace{\mathbf{W}' \begin{pmatrix} \sqrt{\lambda_1} & 0 & \dots & 0 \\ 0 & \sqrt{\lambda_2} & \dots & 0 \\ \vdots & \vdots & \ddots & \vdots \\ 0 & 0 & \dots & \sqrt{\lambda_N} \end{pmatrix}}_{\mathbf{X}} \underbrace{\begin{pmatrix} \sqrt{\lambda_1} & 0 & \dots & 0 \\ 0 & \sqrt{\lambda_2} & \dots & 0 \\ \vdots & \vdots & \ddots & \vdots \\ 0 & 0 & \dots & \sqrt{\lambda_N} \end{pmatrix}}_{\mathbf{X}'} \mathbf{W}\end{aligned}$$

Approximate  $\mathbf{X}$  with first  $L$  eigenvectors

# Classic Multidimensional Scaling $\rightsquigarrow$ Optimization

Define  $\mathbf{z}_i \in \mathbb{R}^L$  as,

$$\mathbf{z}_i^* = \left( \sqrt{\lambda_1} w_{1i}, \sqrt{\lambda_2} w_{2i}, \dots, \sqrt{\lambda_L} w_{Li} \right)$$

# Classic Multidimensional Scaling $\rightsquigarrow$ Optimization

Define  $\mathbf{z}_i \in \mathbb{R}^L$  as,

$$\mathbf{z}_i^* = \left( \sqrt{\lambda_1} w_{1i}, \sqrt{\lambda_2} w_{2i}, \dots, \sqrt{\lambda_L} w_{Li} \right)$$

- Same diagnostics as before: eigenvalues

# Classic Multidimensional Scaling $\rightsquigarrow$ Optimization

Define  $\mathbf{z}_i \in \mathbb{R}^L$  as,

$$\mathbf{z}_i^* = \left( \sqrt{\lambda_1} w_{1i}, \sqrt{\lambda_2} w_{2i}, \dots, \sqrt{\lambda_L} w_{Li} \right)$$

- Same diagnostics as before: eigenvalues
- Same issues in dimension reduction  $\rightsquigarrow$  what are the dimensions for your task

# Classic Multidimensional Scaling $\rightsquigarrow$ Optimization

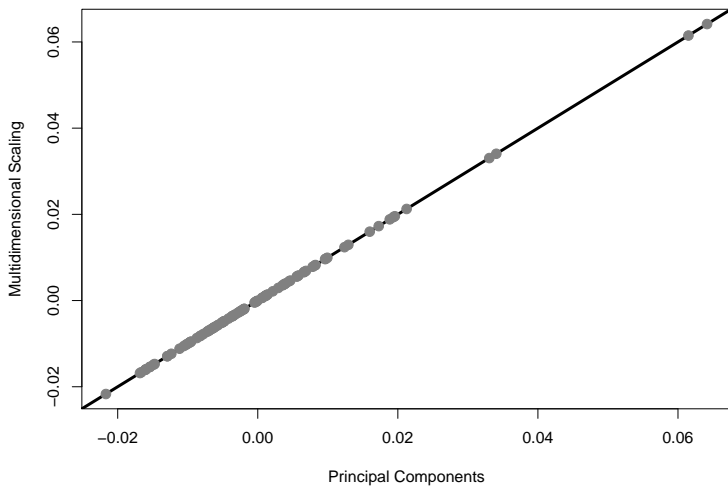
Define  $\mathbf{z}_i \in \mathbb{R}^L$  as,

$$\mathbf{z}_i^* = \left( \sqrt{\lambda_1} w_{1i}, \sqrt{\lambda_2} w_{2i}, \dots, \sqrt{\lambda_L} w_{Li} \right)$$

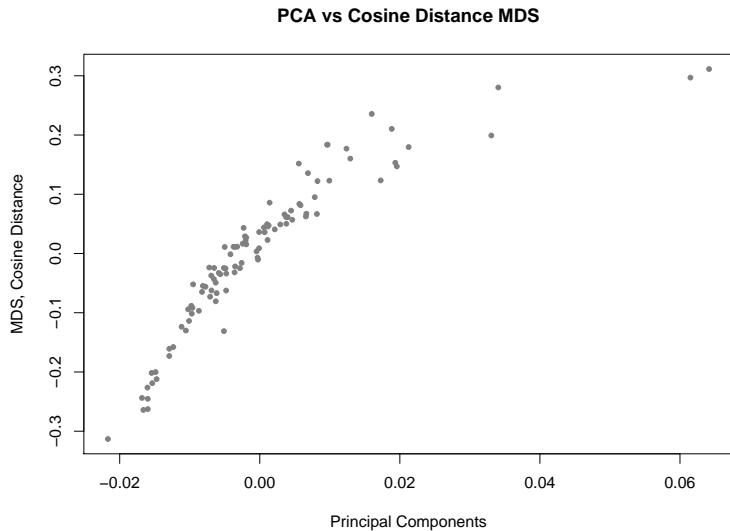
- Same diagnostics as before: eigenvalues
- Same issues in dimension reduction  $\rightsquigarrow$  what are the dimensions for your task
- Information up to **rotation** and arbitrary **center** (same result with any **orthogonal** matrix  $\mathbf{M}$ )

# Comparing MDS and PCA

Comparing PCA and MDS on Senate Data



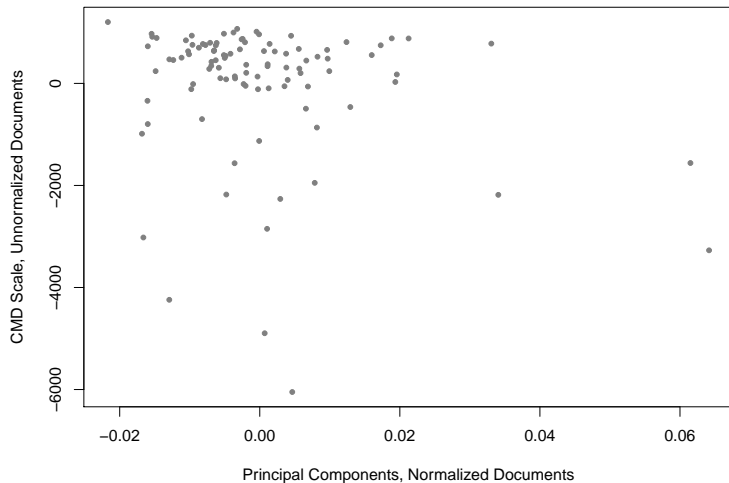
# Comparing MDS and PCA





# Comparing MDS and PCA

PCA vs Euclidean Distance MDS, Unnormalized



# Sammon Multidimensional Scaling $\rightsquigarrow$ Manifold Learning

Many other methods we won't cover:

# Sammon Multidimensional Scaling $\rightsquigarrow$ Manifold Learning

Many other methods we won't cover:

- Sammon Multidimensional Scaling  $\rightsquigarrow$  prioritizes small differences

# Sammon Multidimensional Scaling $\rightsquigarrow$ Manifold Learning

Many other methods we won't cover:

- Sammon Multidimensional Scaling  $\rightsquigarrow$  prioritizes small differences

$$f(\mathbf{X}, \mathbf{Z}) = \sum_{i \neq j} \frac{(d_{ij} - \|z_i - z_j\|)^2}{d_{ij}}$$

# Sammon Multidimensional Scaling $\rightsquigarrow$ Manifold Learning

Many other methods we won't cover:

- Sammon Multidimensional Scaling  $\rightsquigarrow$  prioritizes small differences

$$f(\mathbf{X}, \mathbf{Z}) = \sum_{i \neq j} \frac{(d_{ij} - \|z_i - z_j\|)^2}{d_{ij}}$$

- **Landmark MDS**  $\rightsquigarrow$  for large scale embeddings

# Sammon Multidimensional Scaling $\rightsquigarrow$ Manifold Learning

Many other methods we won't cover:

- Sammon Multidimensional Scaling  $\rightsquigarrow$  prioritizes small differences

$$f(\mathbf{X}, \mathbf{Z}) = \sum_{i \neq j} \frac{(d_{ij} - \|z_i - z_j\|)^2}{d_{ij}}$$

- **Landmark MDS**  $\rightsquigarrow$  for large scale embeddings
- **ISOMap**  $\rightsquigarrow$  Non-linear embedding, emphasizing local patterns

# Sammon Multidimensional Scaling $\rightsquigarrow$ Manifold Learning

Many other methods we won't cover:

- Sammon Multidimensional Scaling  $\rightsquigarrow$  prioritizes small differences

$$f(\mathbf{X}, \mathbf{Z}) = \sum_{i \neq j} \frac{(d_{ij} - \|z_i - z_j\|)^2}{d_{ij}}$$

- **Landmark MDS**  $\rightsquigarrow$  for large scale embeddings
- **ISOMap**  $\rightsquigarrow$  Non-linear embedding, emphasizing local patterns

1) Calculate distance matrix

# Sammon Multidimensional Scaling $\rightsquigarrow$ Manifold Learning

Many other methods we won't cover:

- Sammon Multidimensional Scaling  $\rightsquigarrow$  prioritizes small differences

$$f(\mathbf{X}, \mathbf{Z}) = \sum_{i \neq j} \frac{(d_{ij} - \|z_i - z_j\|)^2}{d_{ij}}$$

- **Landmark MDS**  $\rightsquigarrow$  for large scale embeddings
  - **ISOMap**  $\rightsquigarrow$  Non-linear embedding, emphasizing local patterns
- 1) Calculate distance matrix
  - 2) Create  $K$  nearest neighbor graph: find  $K$  nearest neighbors for each point, with edge weight = Euclidean distance



# Sammon Multidimensional Scaling $\rightsquigarrow$ Manifold Learning

Many other methods we won't cover:

- Sammon Multidimensional Scaling  $\rightsquigarrow$  prioritizes small differences

$$f(\mathbf{X}, \mathbf{Z}) = \sum_{i \neq j} \frac{(d_{ij} - \|z_i - z_j\|)^2}{d_{ij}}$$

- **Landmark MDS**  $\rightsquigarrow$  for large scale embeddings
- **ISOMap**  $\rightsquigarrow$  Non-linear embedding, emphasizing local patterns

- 1) Calculate distance matrix
- 2) Create  $K$  nearest neighbor graph: find  $K$  nearest neighbors for each point, with edge weight = Euclidean distance
- 3) Compute shortest weighted path between each node (distance matrix)

# Sammon Multidimensional Scaling $\rightsquigarrow$ Manifold Learning

Many other methods we won't cover:

- Sammon Multidimensional Scaling  $\rightsquigarrow$  prioritizes small differences

$$f(\mathbf{X}, \mathbf{Z}) = \sum_{i \neq j} \frac{(d_{ij} - \|z_i - z_j\|)^2}{d_{ij}}$$

- **Landmark MDS**  $\rightsquigarrow$  for large scale embeddings
  - **ISOMap**  $\rightsquigarrow$  Non-linear embedding, emphasizing local patterns
- 1) Calculate distance matrix
  - 2) Create  $K$  nearest neighbor graph: find  $K$  nearest neighbors for each point, with edge weight = Euclidean distance
  - 3) Compute shortest weighted path between each node (distance matrix)
  - 4) Apply classic MDS

# Sammon Multidimensional Scaling $\rightsquigarrow$ Manifold Learning

Many other methods we won't cover:

- Sammon Multidimensional Scaling  $\rightsquigarrow$  prioritizes small differences

$$f(\mathbf{X}, \mathbf{Z}) = \sum_{i \neq j} \frac{(d_{ij} - \|z_i - z_j\|)^2}{d_{ij}}$$

- **Landmark MDS**  $\rightsquigarrow$  for large scale embeddings
- **ISOMap**  $\rightsquigarrow$  Non-linear embedding, emphasizing local patterns

- 1) Calculate distance matrix
- 2) Create  $K$  nearest neighbor graph: find  $K$  nearest neighbors for each point, with edge weight = Euclidean distance
- 3) Compute shortest weighted path between each node (distance matrix)
- 4) Apply classic MDS

Common theme  $\rightsquigarrow$  Manifold Learning

# Sammon Multidimensional Scaling $\rightsquigarrow$ Manifold Learning

Many other methods we won't cover:

- Sammon Multidimensional Scaling  $\rightsquigarrow$  prioritizes small differences

$$f(\mathbf{X}, \mathbf{Z}) = \sum_{i \neq j} \frac{(d_{ij} - \|z_i - z_j\|)^2}{d_{ij}}$$

- **Landmark MDS**  $\rightsquigarrow$  for large scale embeddings
- **ISOMap**  $\rightsquigarrow$  Non-linear embedding, emphasizing local patterns

- 1) Calculate distance matrix
- 2) Create  $K$  nearest neighbor graph: find  $K$  nearest neighbors for each point, with edge weight = Euclidean distance
- 3) Compute shortest weighted path between each node (distance matrix)
- 4) Apply classic MDS

Common theme  $\rightsquigarrow$  Manifold Learning

# Low Dimensional Embedding

- 1) Find lower dimensional space to summarize documents  $\rightsquigarrow$  Eigenvectors
- 2) Thursday: basic language models and the Dirichlet distribution