

Machine Learning 2

DS 4420 - Spring 2020

Dimensionality reduction I

Byron C Wallace



Machine Learning 2

DS 4420 - Spring 2020



Some slides today borrowing from:
Percy Liang (Stanford)

Other material from the MML book (Faisal and Ong)

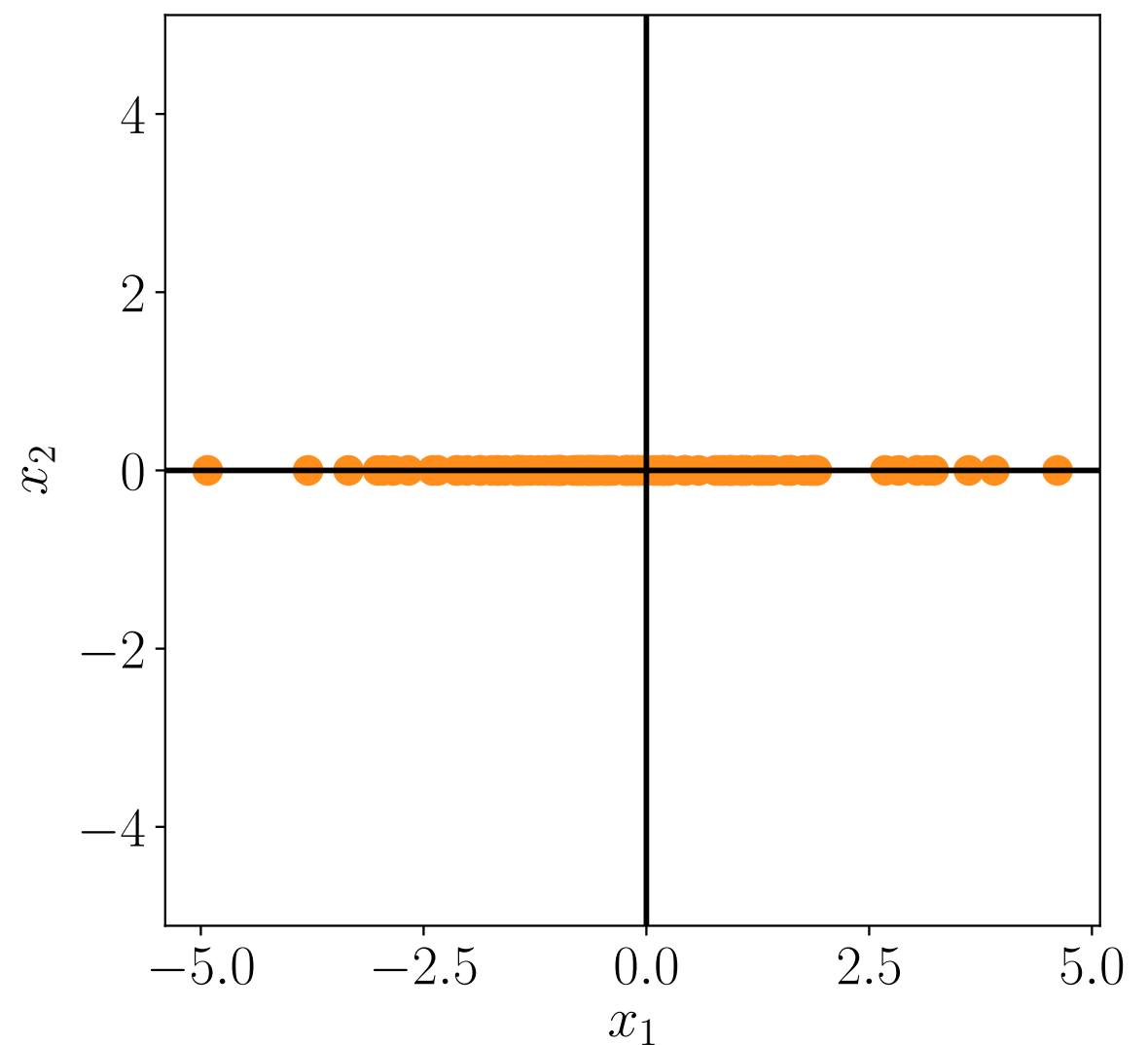
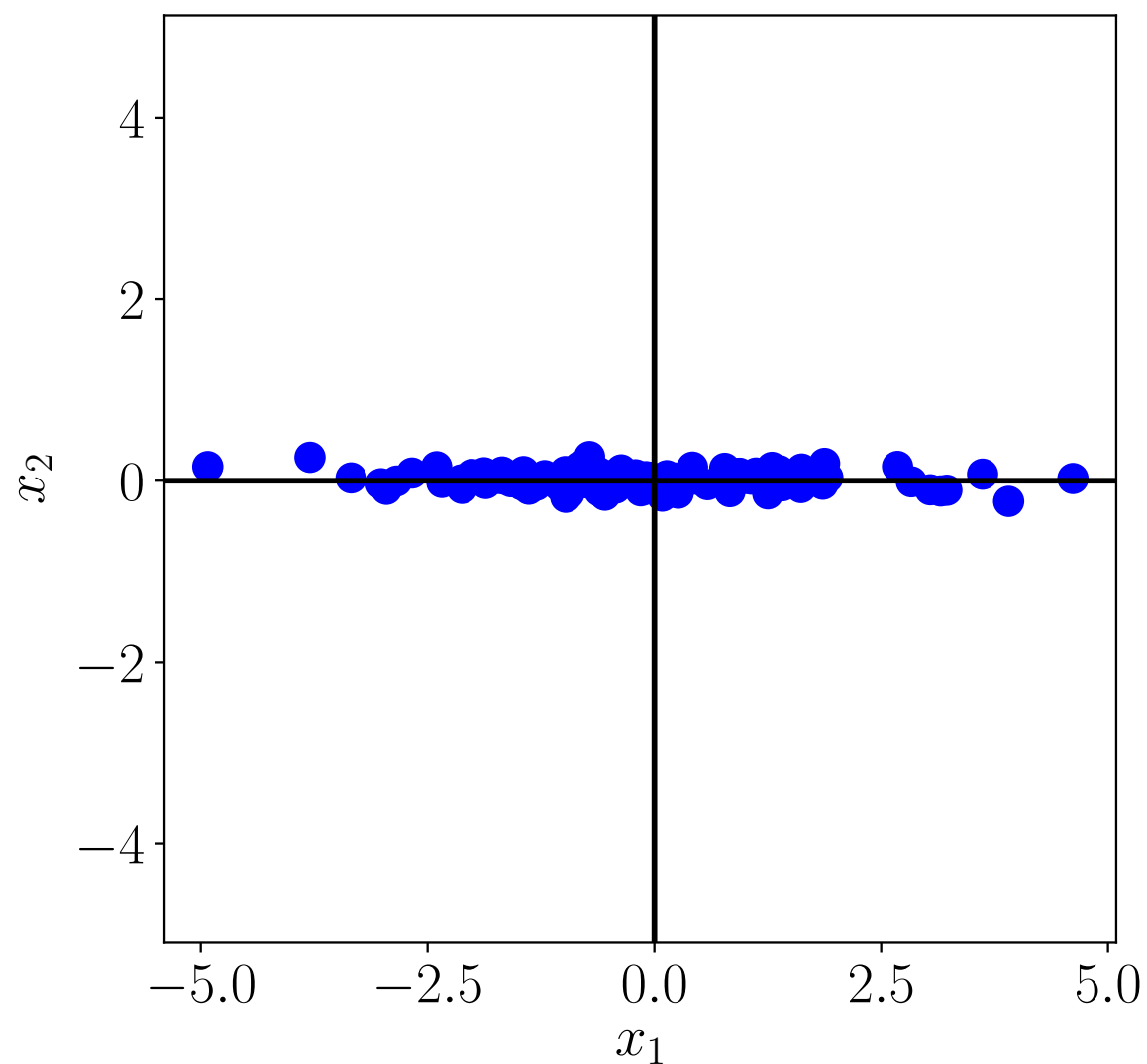


Motivation

- We often want to work with *high dimensional* data (e.g., images). We also often have lots of it.
- This is computationally expensive to store and work with.

Dimensionality Reduction

Fundamental idea Exploit *redundancy* in the data; find *lower-dimensional* representation

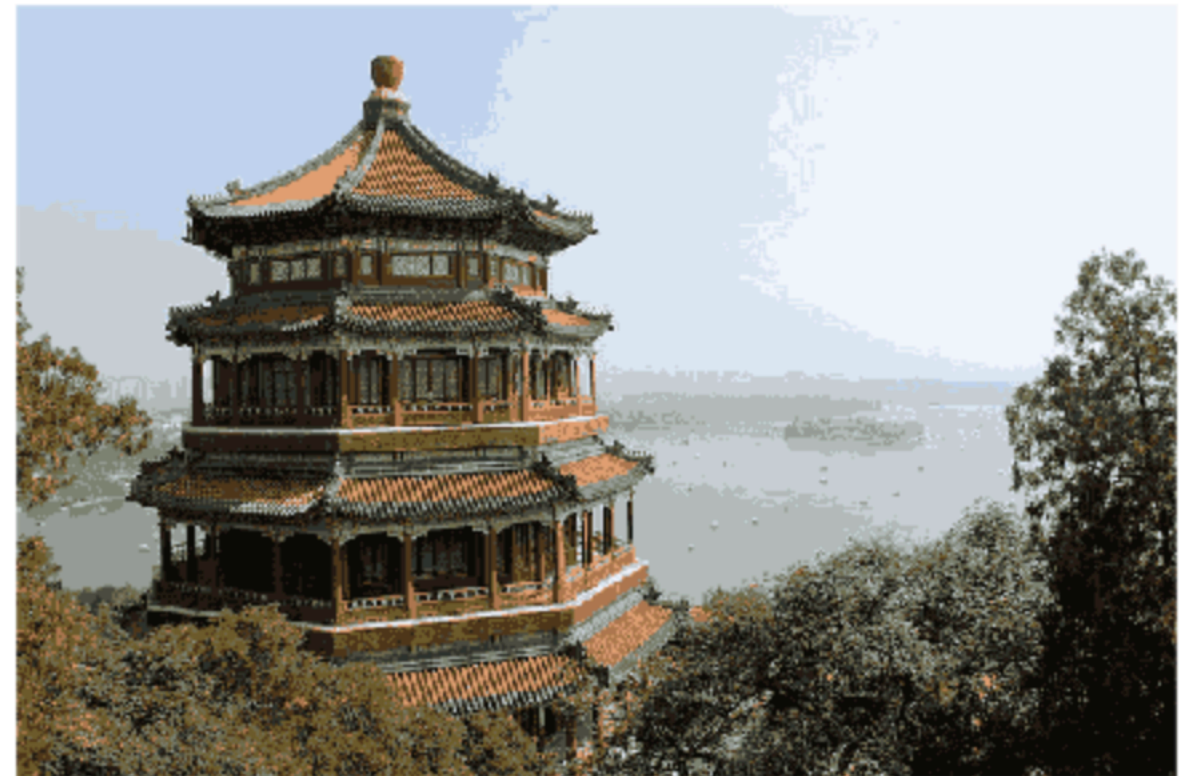


Example (from lecture 5): Dimensionality reduction via k -means

Original Image



16-color Image

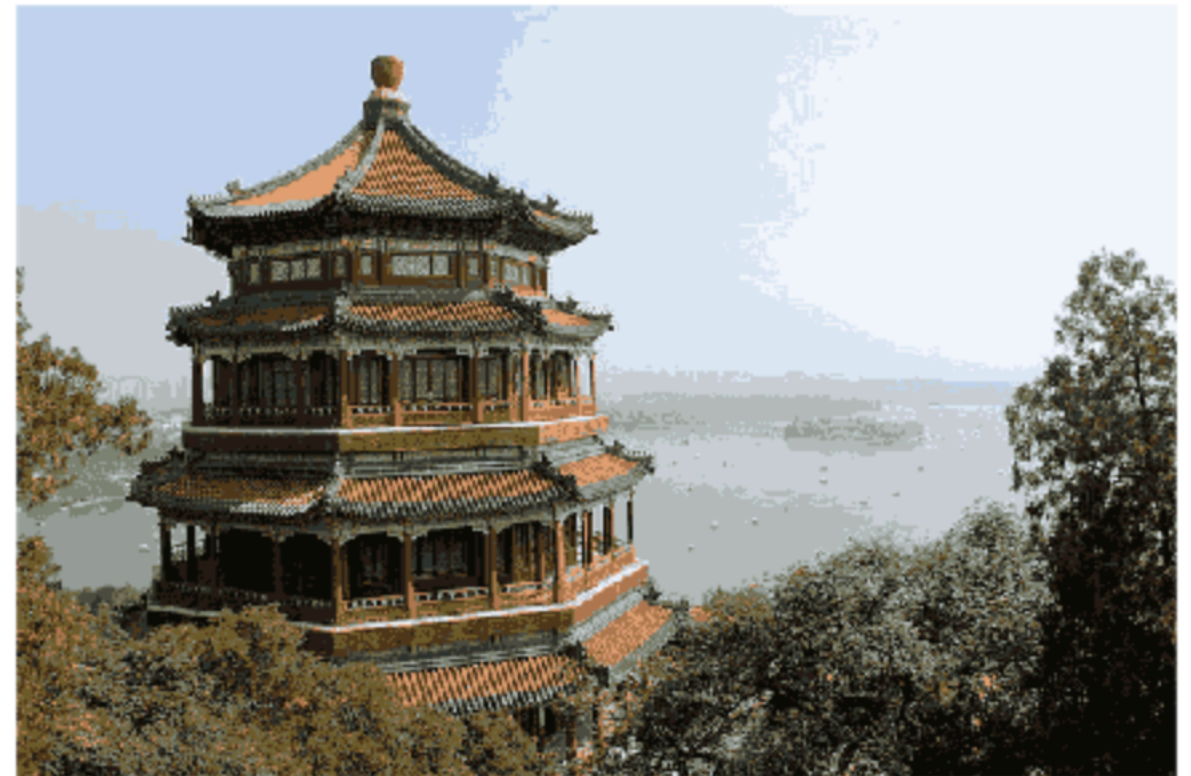


Example (from lecture 5): Dimensionality reduction via k -means

Original Image



16-color Image

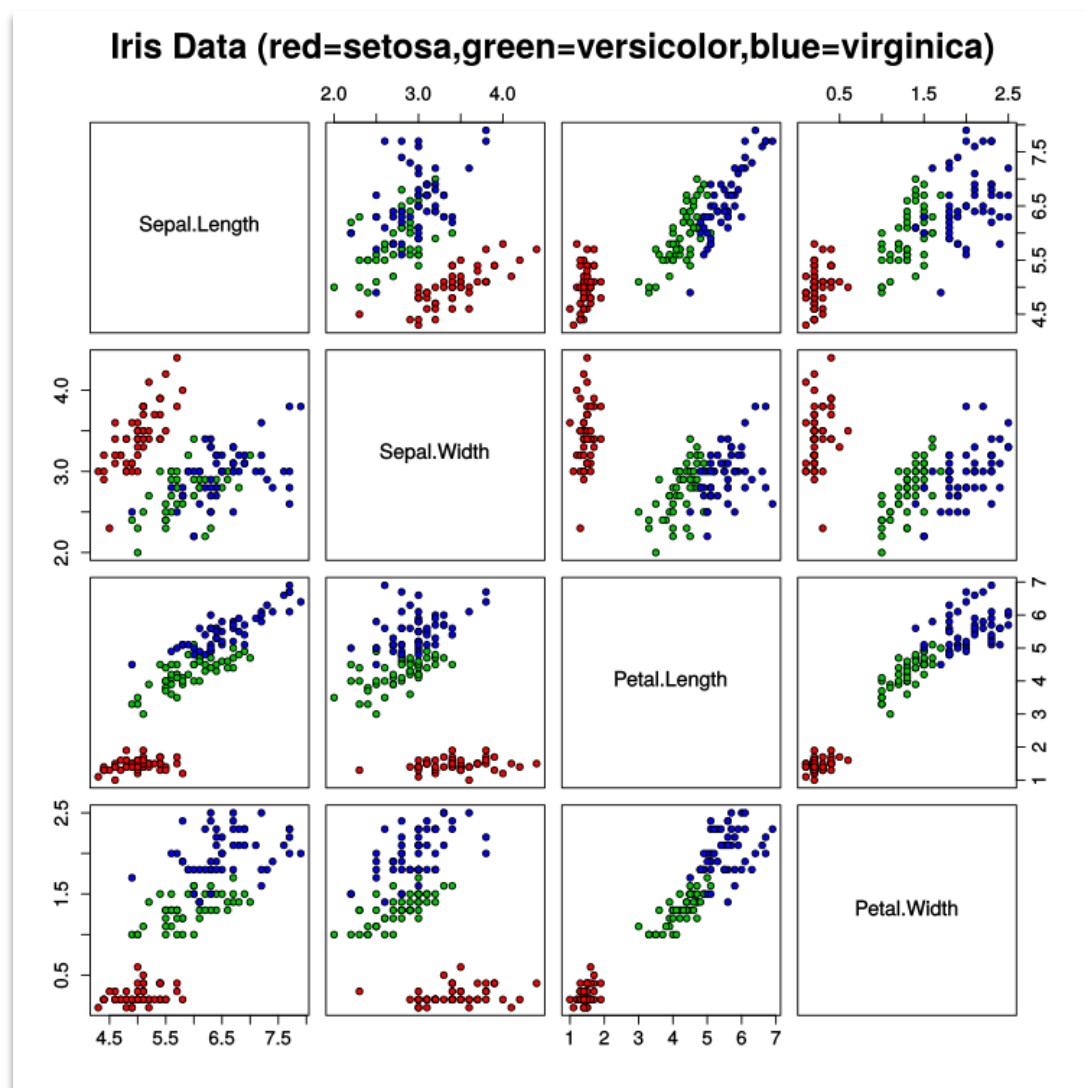


This highlights the natural connection between dimensionality reduction and *compression*.

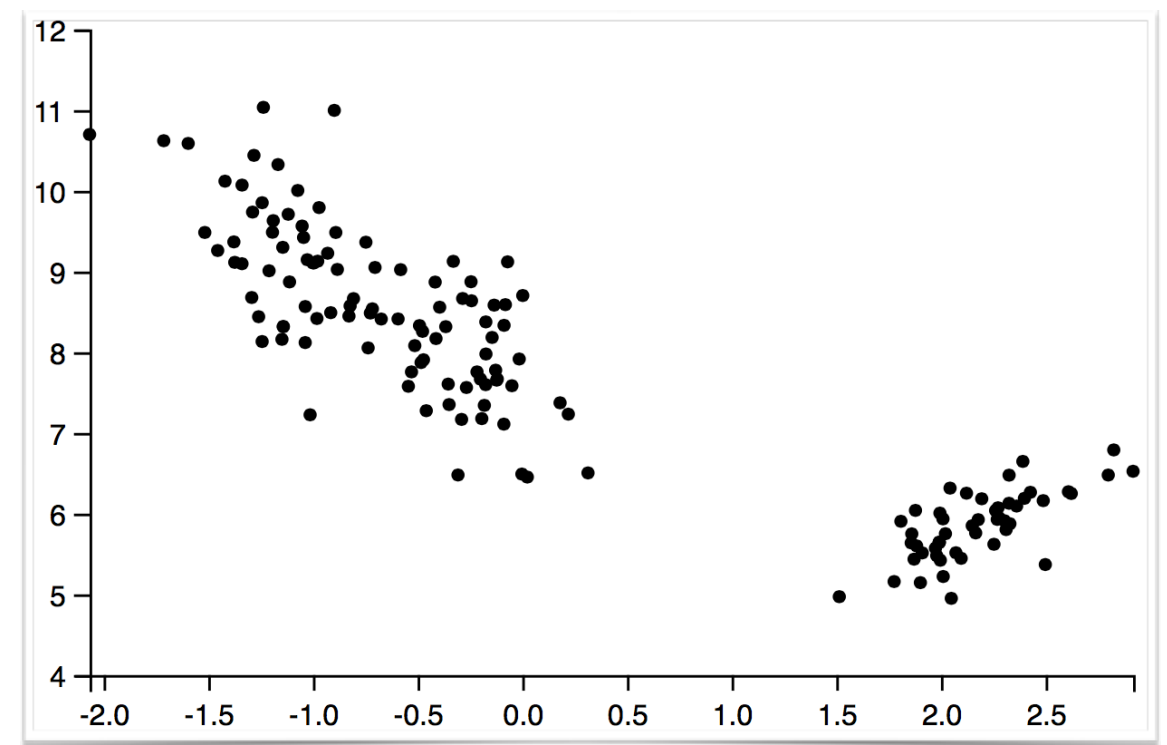
Dimensionality reduction

Goal: Map high dimensional data onto lower-dimensional data in a manner that preserves *distances/similarities*

Original Data (4 dims)



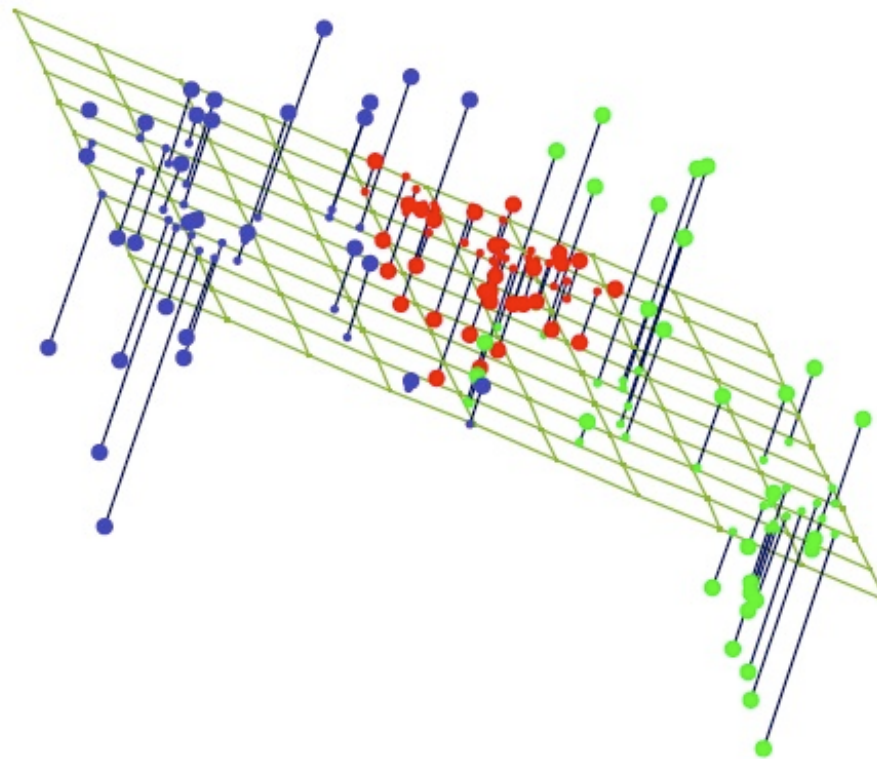
Projection with PCA (2 dims)



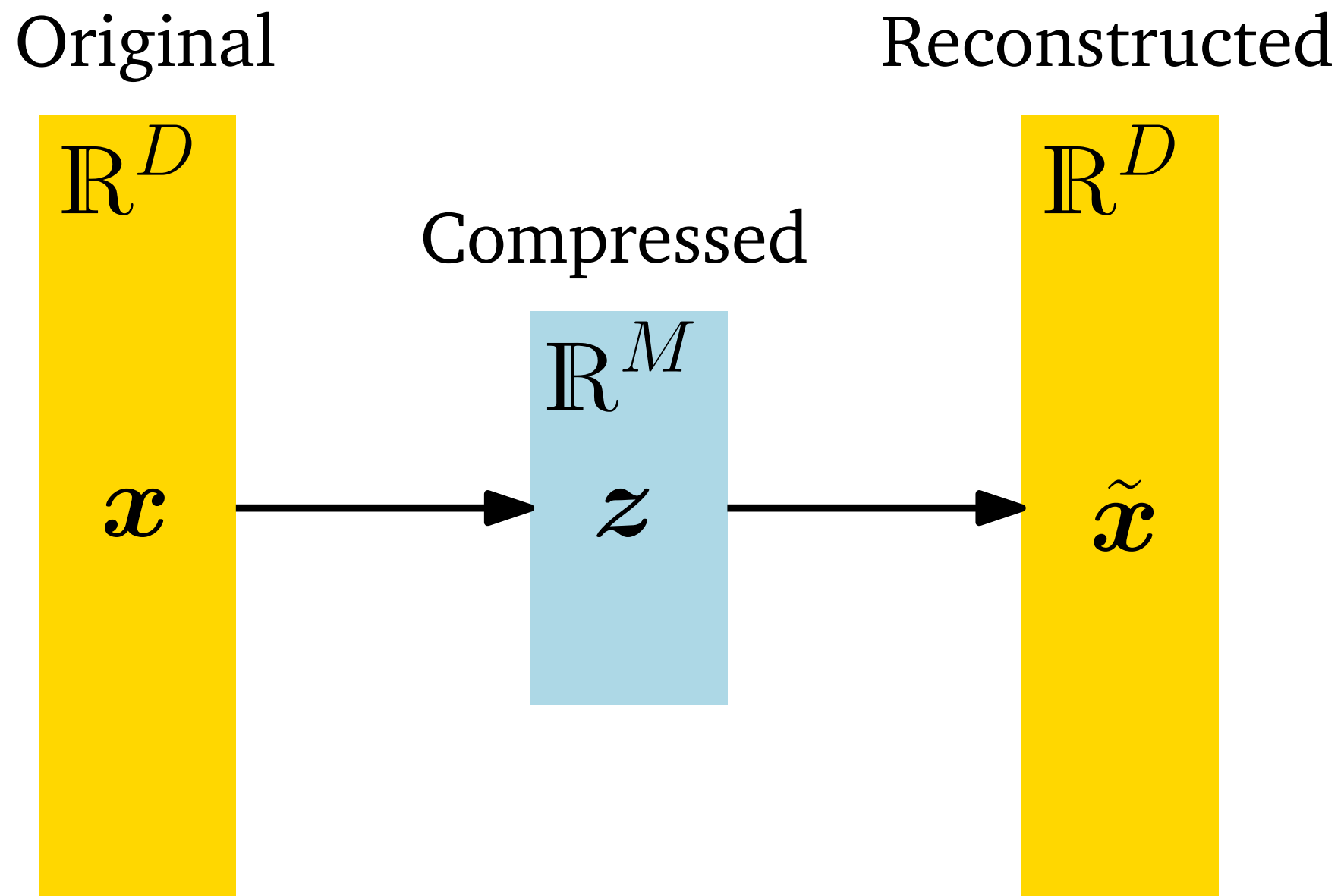
Objective: projection should “preserve” relative distances

Linear dimensionality reduction

Idea: Project high-dimensional vector onto a lower dimensional space



Linear dimensionality reduction



Objective

Key intuition:

$$\underbrace{\text{variance of data}}_{\text{fixed}} = \underbrace{\text{captured variance}}_{\text{want large}} + \underbrace{\text{reconstruction error}}_{\text{want small}}$$

Principal Component Analysis (on board)

In Sum: Principal Component Analysis

Data

$$\mathbf{X} = \begin{pmatrix} | & & | \\ \mathbf{x}_1 & \cdots & \mathbf{x}_n \\ | & & | \end{pmatrix} \in \mathbb{R}^{d \times n}$$

Orthonormal Basis

$$\mathbf{U} = \begin{pmatrix} | & & | \\ \mathbf{u}_1 & \cdots & \mathbf{u}_d \\ | & & | \end{pmatrix} \in \mathbb{R}^{d \times d}$$

Eigenvectors of Covariance

$$\mathbf{C} = \frac{1}{n} \sum_{j=1}^n \mathbf{x}_j \mathbf{x}_j^\top = \frac{1}{n} \mathbf{X} \mathbf{X}^\top$$

$$\mathbf{C} \mathbf{u}_j = \lambda_j \mathbf{u}_j$$

Eigen-decomposition

$$\mathbf{C} = \mathbf{U} \mathbf{\Lambda} \mathbf{U}^\top$$

$$\mathbf{\Lambda} = \begin{pmatrix} \lambda_1 & & & \\ & \lambda_2 & & \\ & & \cdots & \\ & & & \lambda_d \end{pmatrix}$$

Idea: Take **top- k** eigenvectors to maximize variance

Getting the eigenvalues, two ways

- Direct eigenvalue decomposition of the covariance matrix

$$\mathbf{S} = \frac{1}{N} \sum_{n=1}^N \mathbf{x}_n \mathbf{x}_n^\top = \frac{1}{N} \mathbf{X} \mathbf{X}^\top$$

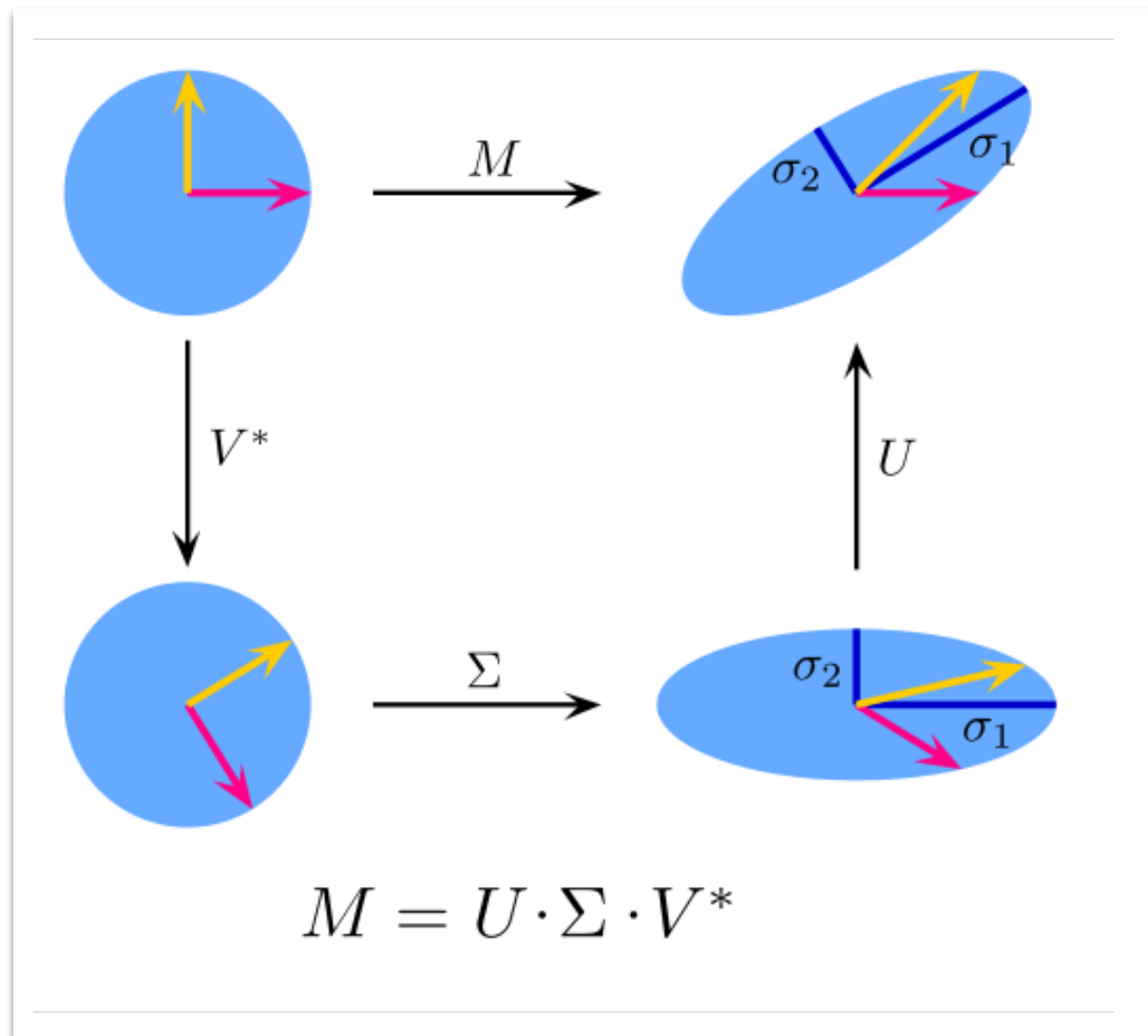
Getting the eigenvalues, two ways

- Direct eigenvalue decomposition of the covariance matrix

$$\mathbf{S} = \frac{1}{N} \sum_{n=1}^N \mathbf{x}_n \mathbf{x}_n^\top = \frac{1}{N} \mathbf{X} \mathbf{X}^\top$$

- Singular Value Decomposition (SVD)

Singular Value Decomposition



Idea: Decompose the $d \times n$ matrix \mathbf{X} into

1. A $n \times n$ basis \mathbf{V} (unitary matrix)
2. A $d \times n$ matrix Σ (diagonal projection)
3. A $d \times d$ basis \mathbf{U} (unitary matrix)

$$\mathbf{X} = \mathbf{U}_{d \times d} \Sigma_{d \times n} \mathbf{V}_{n \times n}^T$$

SVD for PCA

$$\underbrace{\mathbf{X}}_{D \times N} = \underbrace{\mathbf{U}}_{D \times D} \underbrace{\mathbf{\Sigma}}_{D \times N} \underbrace{\mathbf{V}^\top}_{N \times N}$$

$$\mathbf{S} = \frac{1}{N} \mathbf{X} \mathbf{X}^\top = \frac{1}{N} \mathbf{U} \mathbf{\Sigma} \underbrace{\mathbf{V}^\top \mathbf{V}}_{=\mathbf{I}_N} \mathbf{\Sigma}^\top \mathbf{U}^\top = \frac{1}{N} \mathbf{U} \mathbf{\Sigma} \mathbf{\Sigma}^\top \mathbf{U}^\top$$

SVD for PCA

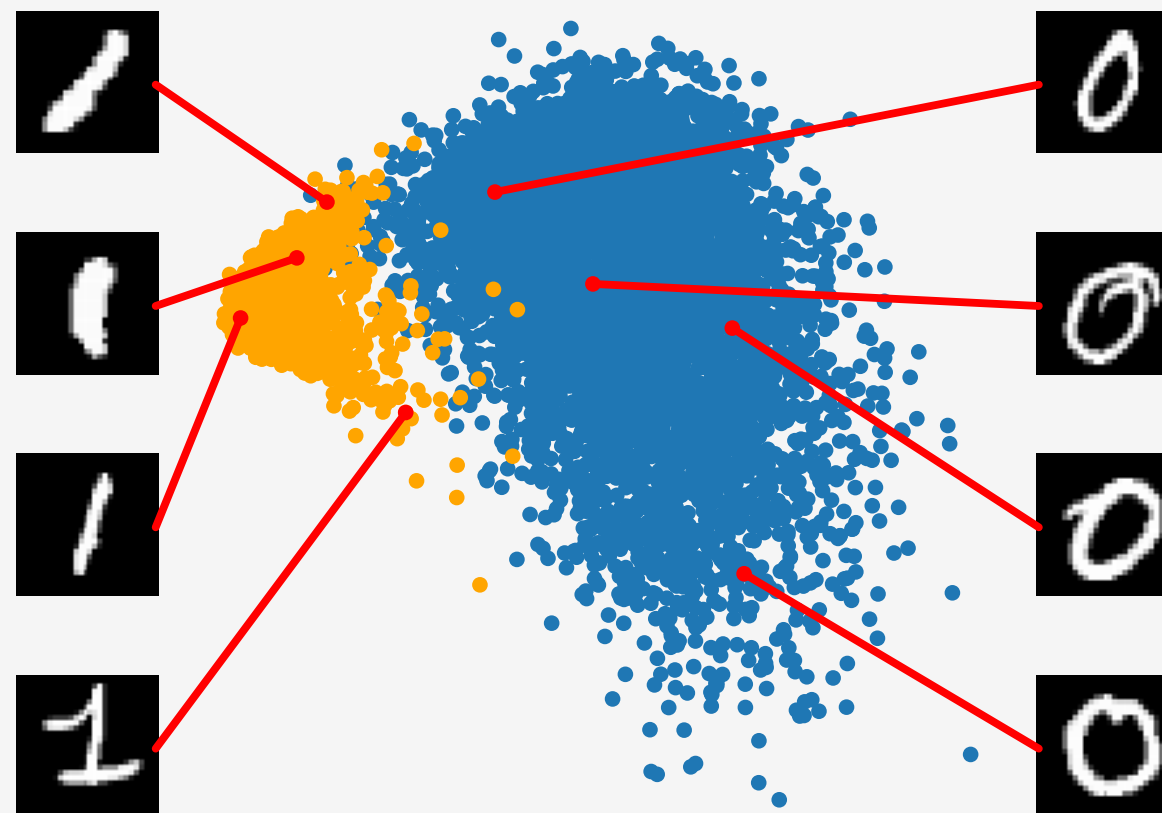
$$\underbrace{\mathbf{X}}_{D \times N} = \underbrace{\mathbf{U}}_{D \times D} \underbrace{\mathbf{\Sigma}}_{D \times N} \underbrace{\mathbf{V}^\top}_{N \times N}$$

$$\mathbf{S} = \frac{1}{N} \mathbf{X} \mathbf{X}^\top = \frac{1}{N} \mathbf{U} \mathbf{\Sigma} \underbrace{\mathbf{V}^\top \mathbf{V}}_{=\mathbf{I}_N} \mathbf{\Sigma}^\top \mathbf{U}^\top = \frac{1}{N} \mathbf{U} \mathbf{\Sigma} \mathbf{\Sigma}^\top \mathbf{U}^\top$$

It turns out the columns of \mathbf{U} are the eigenvectors of $\mathbf{X} \mathbf{X}^\top$

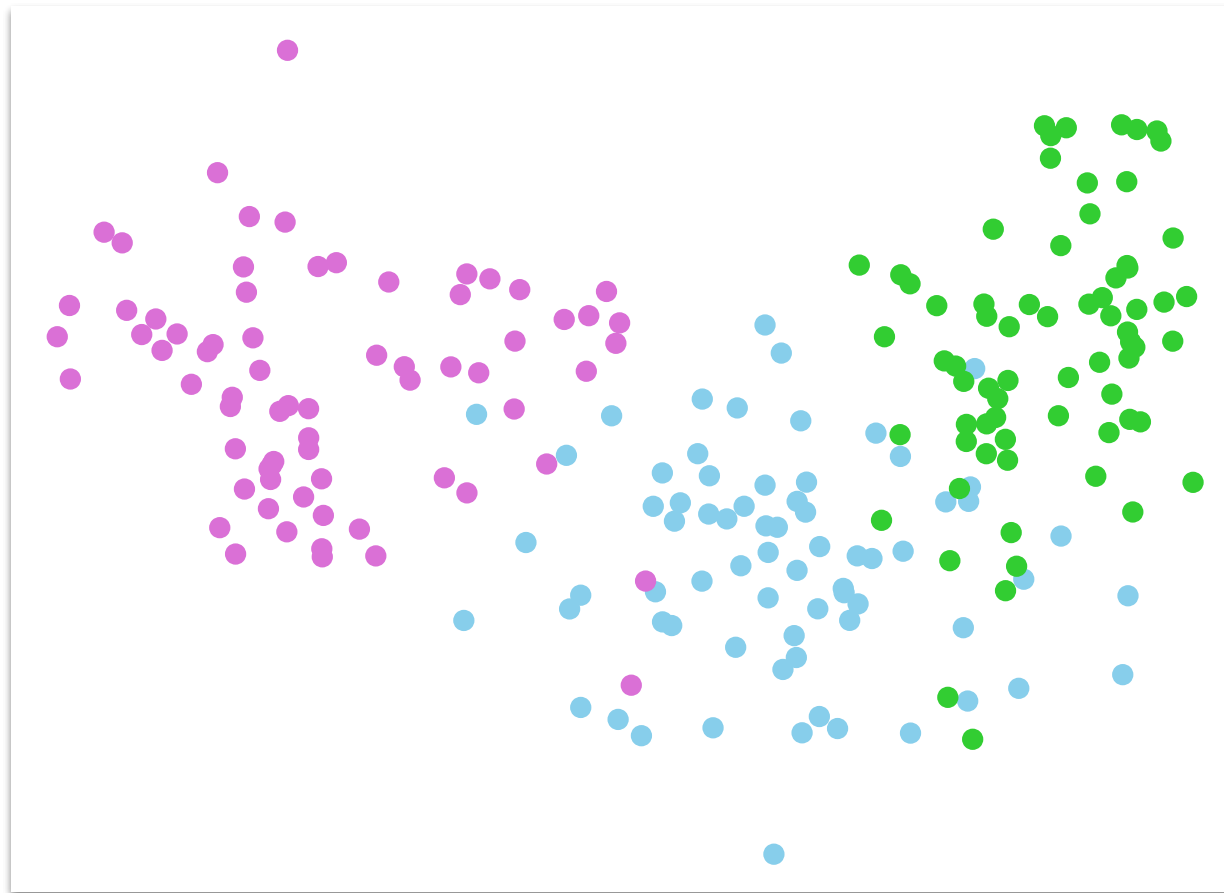
Principal Component Analysis

Example 10.3 (MNIST Digits Embedding)

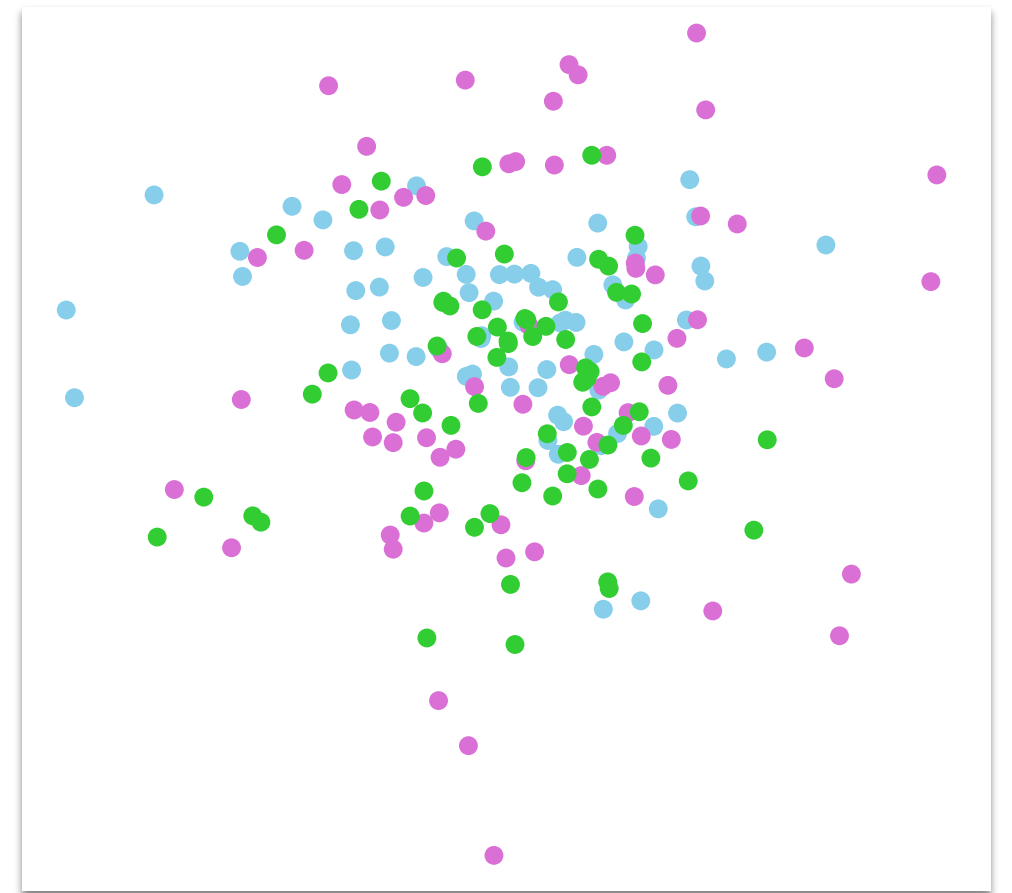


Principal Component Analysis

Top 2 components



Bottom 2 components



Data: three varieties of wheat: Kama, Rosa, Canadian

Attributes: Area, Perimeter, Compactness, Length of Kernel, Width of Kernel, Asymmetry Coefficient, Length of Groove

Eigen-faces [Turk & Pentland 1991]

- d = number of pixels
- Each $\mathbf{x}_i \in \mathbb{R}^d$ is a face image
- x_{ji} = intensity of the j -th pixel in image i

Eigen-faces [Turk & Pentland 1991]

- d = number of pixels
- Each $\mathbf{x}_i \in \mathbb{R}^d$ is a face image
- \mathbf{x}_{ji} = intensity of the j -th pixel in image i

$$\begin{matrix} \mathbf{X}_{d \times n} & \approx & \mathbf{U}_{d \times k} & \mathbf{Z}_{k \times n} \\ \left(\begin{array}{c} \text{[Image 1]} \quad \dots \quad \text{[Image } n\text{]} \end{array} \right) & \approx & \left(\begin{array}{c} \text{[Eigenface 1]} \quad \text{[Eigenface 2]} \quad \dots \quad \text{[Eigenface } k\text{]} \end{array} \right) & \left(\begin{array}{c} \mathbf{z}_1 \quad \dots \quad \mathbf{z}_n \end{array} \right) \end{matrix}$$

Eigen-faces [Turk & Pentland 1991]

- d = number of pixels
- Each $\mathbf{x}_i \in \mathbb{R}^d$ is a face image
- \mathbf{x}_{ji} = intensity of the j -th pixel in image i

$$\begin{matrix} \mathbf{X}_{d \times n} & \approx & \mathbf{U}_{d \times k} & \mathbf{Z}_{k \times n} \\ \left(\begin{array}{c} \text{[Image 1]} \quad \dots \quad \text{[Image n]} \end{array} \right) & \approx & \left(\begin{array}{c} \text{[Eigenface 1]} \quad \text{[Eigenface 2]} \quad \dots \quad \text{[Eigenface k]} \end{array} \right) & \left(\begin{array}{c} \mathbf{z}_1 \quad \dots \quad \mathbf{z}_n \end{array} \right) \end{matrix}$$

Idea: \mathbf{z}_i more “meaningful” representation of i -th face than \mathbf{x}_i

Can use \mathbf{z}_i for nearest-neighbor classification

Eigen-faces [Turk & Pentland 1991]

- d = number of pixels
- Each $\mathbf{x}_i \in \mathbb{R}^d$ is a face image
- \mathbf{x}_{ji} = intensity of the j -th pixel in image i

$$\begin{array}{c} \mathbf{X}_{d \times n} \\ \left(\begin{array}{c} \text{[Face 1]} \quad \dots \quad \text{[Face } n \text{]} \end{array} \right) \end{array} \approx \begin{array}{c} \mathbf{U}_{d \times k} \\ \left(\begin{array}{c} \text{[Eigenface 1]} \quad \text{[Eigenface 2]} \quad \dots \quad \text{[Eigenface } k \text{]} \end{array} \right) \end{array} \begin{array}{c} \mathbf{Z}_{k \times n} \\ \left(\begin{array}{c} | \quad | \quad \dots \quad | \\ \mathbf{z}_1 \quad \dots \quad \mathbf{z}_n \\ | \quad | \quad \dots \quad | \end{array} \right) \end{array}$$

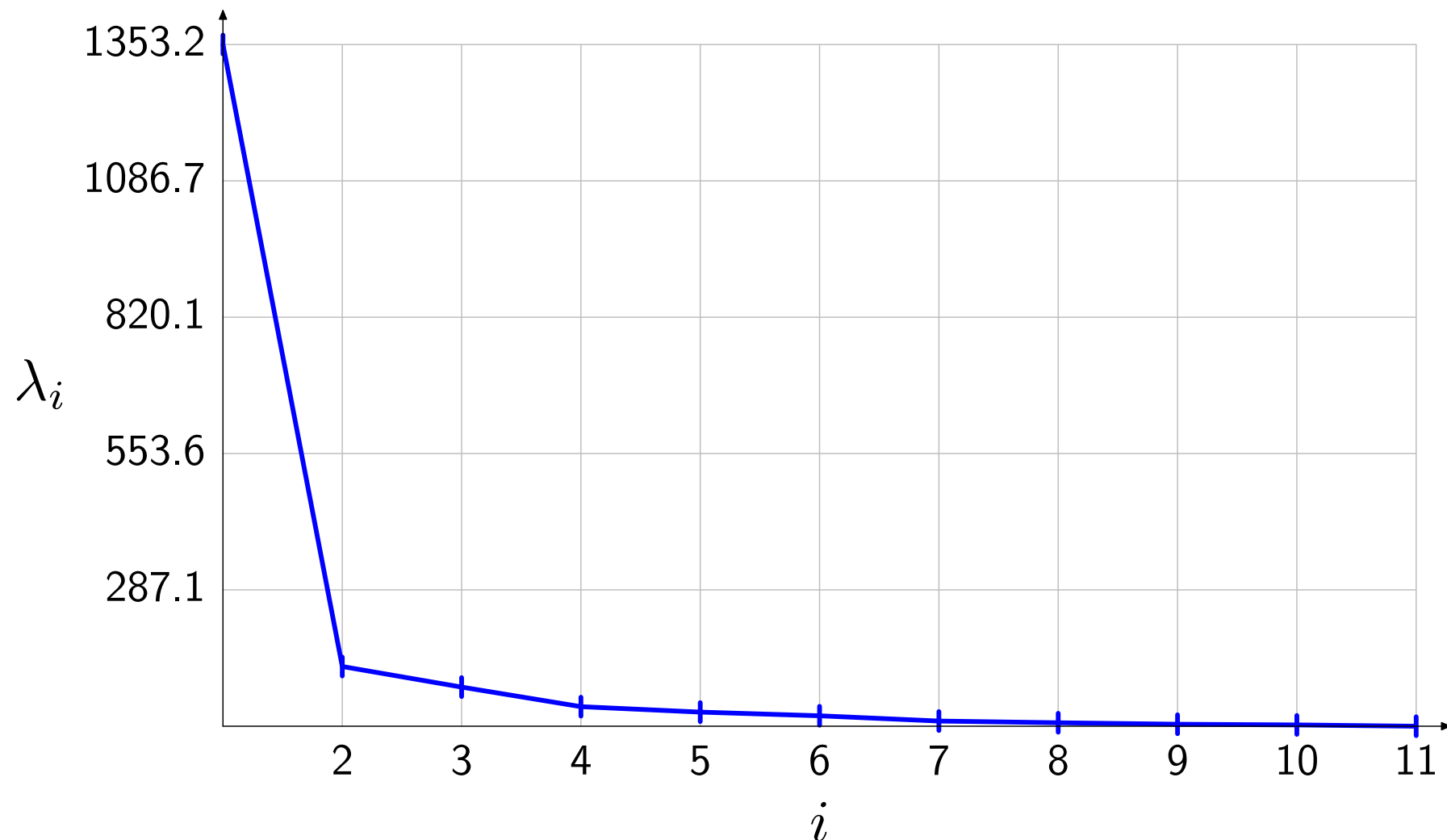
Idea: \mathbf{z}_i more “meaningful” representation of i -th face than \mathbf{x}_i

Can use \mathbf{z}_i for nearest-neighbor classification

Much faster: $O(dk + nk)$ time instead of $O(dn)$ when $n, d \gg k$

Aside: How many components?

- Magnitude of eigenvalues indicate fraction of variance captured.
- Eigenvalues on a face image dataset:



- Eigenvalues typically drop off sharply, so don't need that many.
- Of course variance isn't everything...

Latent Semantic Analysis [Deerwater 1990]

- d = number of words in the vocabulary
- Each $\mathbf{x}_i \in \mathbb{R}^d$ is a vector of word counts

Latent Semantic Analysis [Deerwater 1990]

- d = number of words in the vocabulary
- Each $\mathbf{x}_i \in \mathbb{R}^d$ is a vector of word counts
- \mathbf{x}_{ji} = frequency of word j in document i

$$\begin{array}{c}
 \mathbf{X}_{d \times n} \\
 \left(\begin{array}{cccc}
 \text{stocks: } 2 & \dots & \dots & 0 \\
 \text{chairman: } 4 & \dots & \dots & 1 \\
 \text{the: } 8 & \dots & \dots & 7 \\
 \dots & \vdots & \dots & \vdots \\
 \text{wins: } 0 & \dots & \dots & 2 \\
 \text{game: } 1 & \dots & \dots & 3
 \end{array} \right)
 \end{array}
 \approx
 \begin{array}{c}
 \mathbf{U}_{d \times k} \\
 \left(\begin{array}{cc}
 0.4 & \dots & -0.001 \\
 0.8 & \dots & 0.03 \\
 0.01 & \dots & 0.04 \\
 \vdots & \dots & \vdots \\
 0.002 & \dots & 2.3 \\
 0.003 & \dots & 1.9
 \end{array} \right)
 \end{array}
 \begin{array}{c}
 \mathbf{Z}_{k \times n} \\
 \left(\begin{array}{ccc}
 | & & | \\
 \mathbf{z}_1 & \dots & \mathbf{z}_n \\
 | & & |
 \end{array} \right)
 \end{array}$$

Latent Semantic Analysis [Deerwater 1990]

- d = number of words in the vocabulary
- Each $\mathbf{x}_i \in \mathbb{R}^d$ is a vector of word counts
- \mathbf{x}_{ji} = frequency of word j in document i

$$\begin{array}{c}
 \mathbf{X}_{d \times n} \\
 \left(\begin{array}{cccc}
 \text{stocks: } 2 & \dots & \dots & 0 \\
 \text{chairman: } 4 & \dots & \dots & 1 \\
 \text{the: } 8 & \dots & \dots & 7 \\
 \dots & \vdots & \dots & \vdots \\
 \text{wins: } 0 & \dots & \dots & 2 \\
 \text{game: } 1 & \dots & \dots & 3
 \end{array} \right)
 \end{array}
 \approx
 \begin{array}{c}
 \mathbf{U}_{d \times k} \\
 \left(\begin{array}{cc}
 0.4 & \dots & -0.001 \\
 0.8 & \dots & 0.03 \\
 0.01 & \dots & 0.04 \\
 \vdots & \dots & \vdots \\
 0.002 & \dots & 2.3 \\
 0.003 & \dots & 1.9
 \end{array} \right)
 \end{array}
 \begin{array}{c}
 \mathbf{Z}_{k \times n} \\
 \left(\begin{array}{ccc}
 | & & | \\
 \mathbf{z}_1 & \dots & \mathbf{z}_n \\
 | & & |
 \end{array} \right)
 \end{array}$$

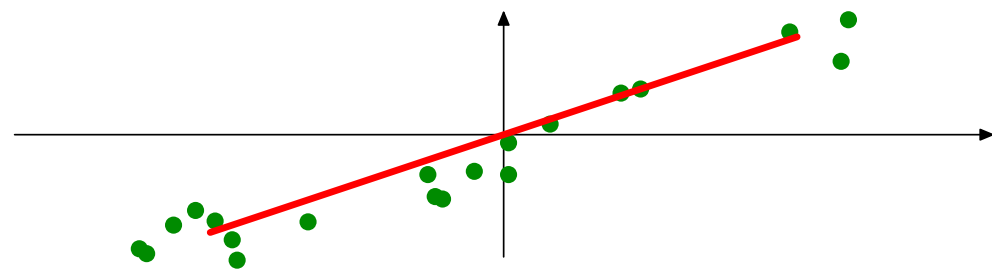
How to measure similarity between two documents?

$$\mathbf{z}_1^\top \mathbf{z}_2 \text{ is probably better than } \mathbf{x}_1^\top \mathbf{x}_2$$

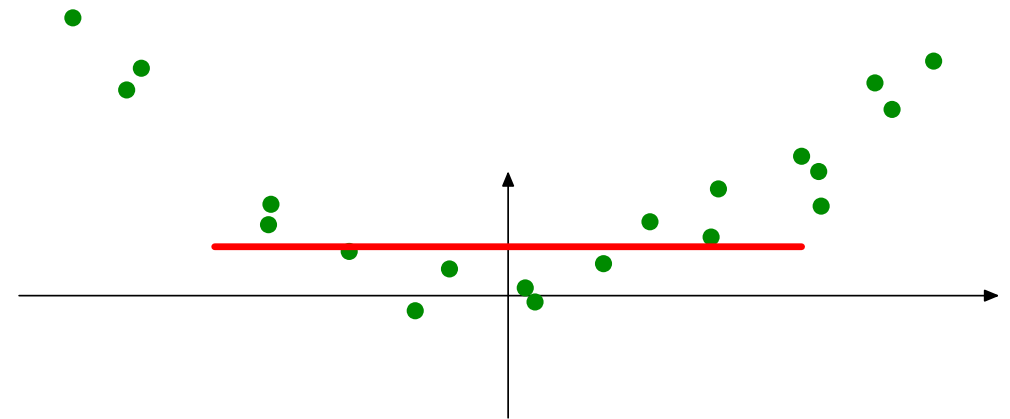
Probabilistic PCA

- If we define a *prior* over z then we can **sample** from the latent space and hallucinate images

Limitations of Linearity

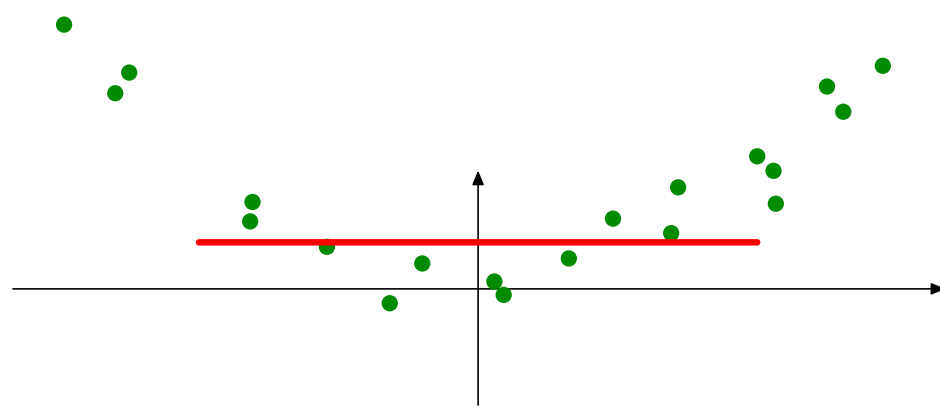


PCA is effective

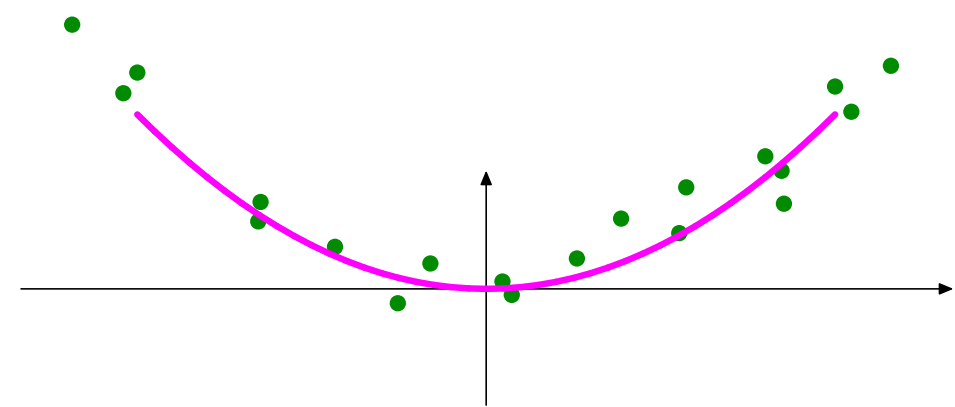


PCA is ineffective

Nonlinear PCA

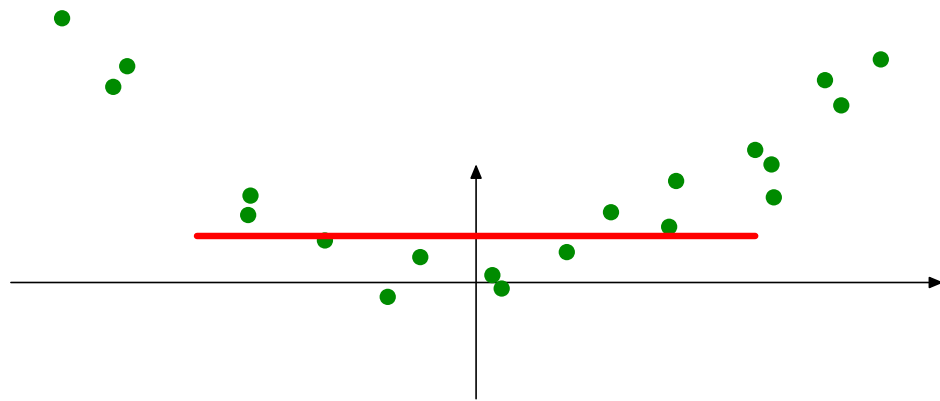


Broken solution

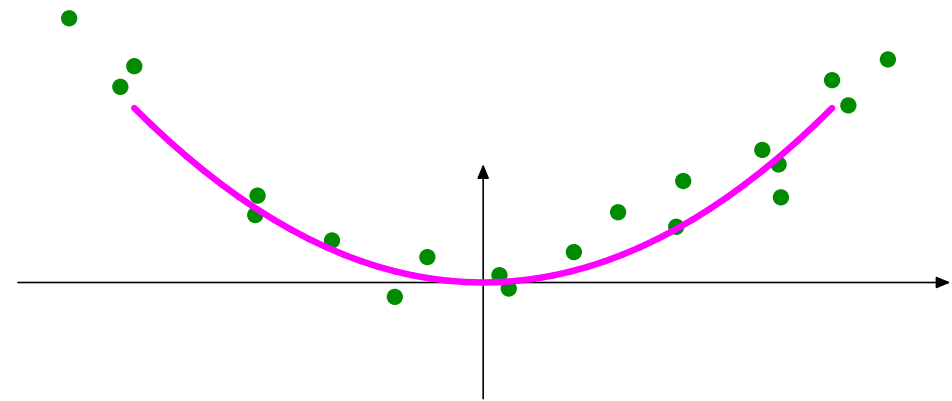


Desired solution

Nonlinear PCA



Broken solution



Desired solution

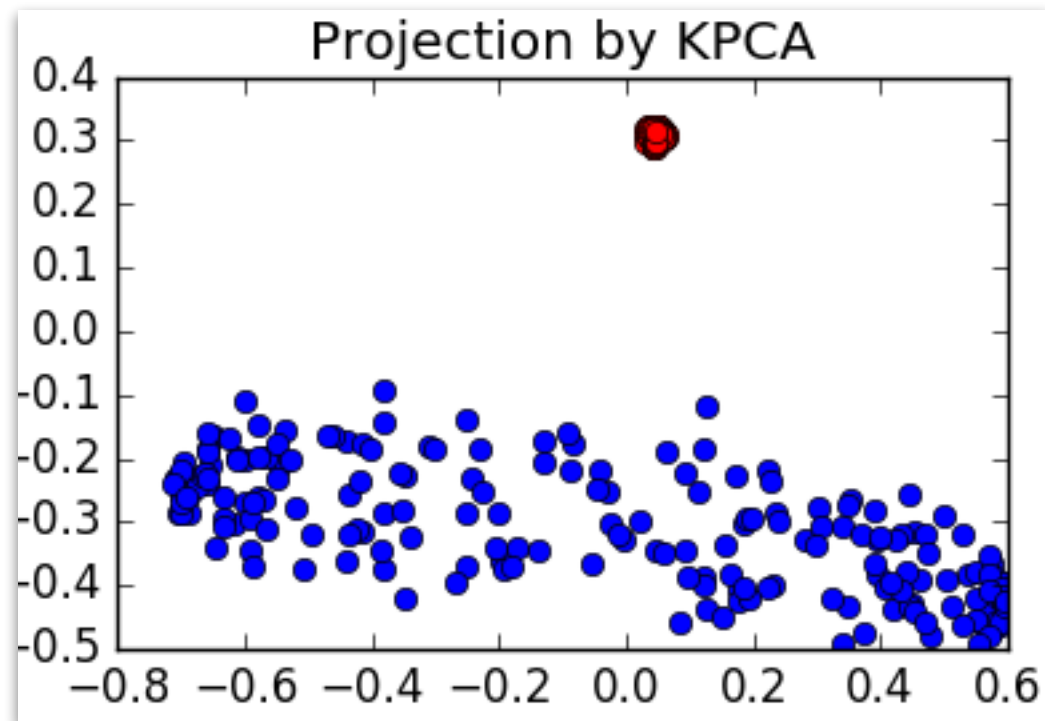
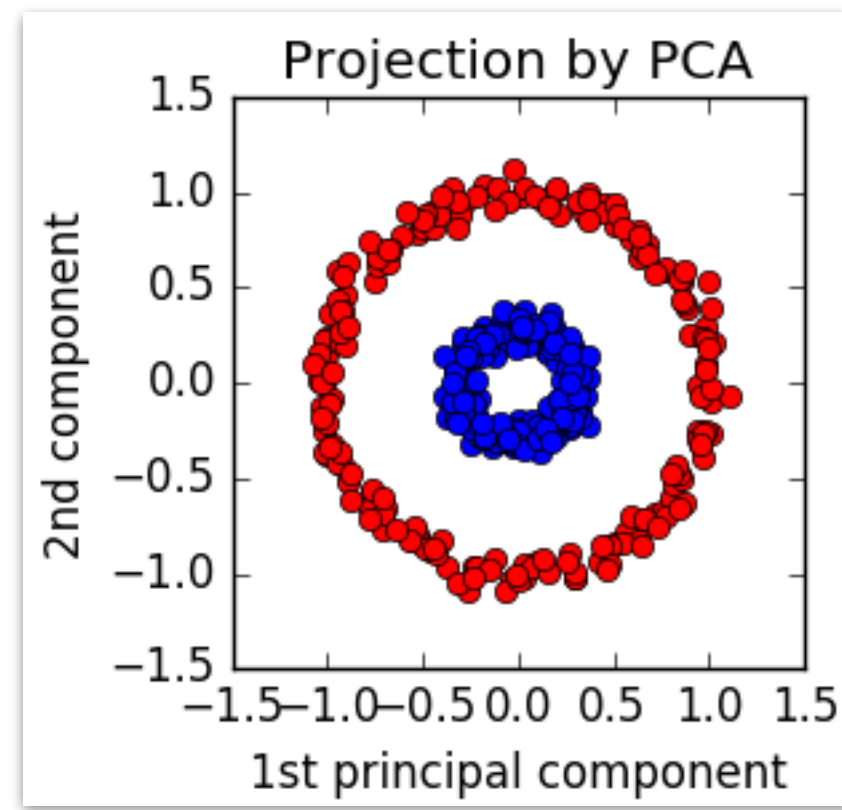
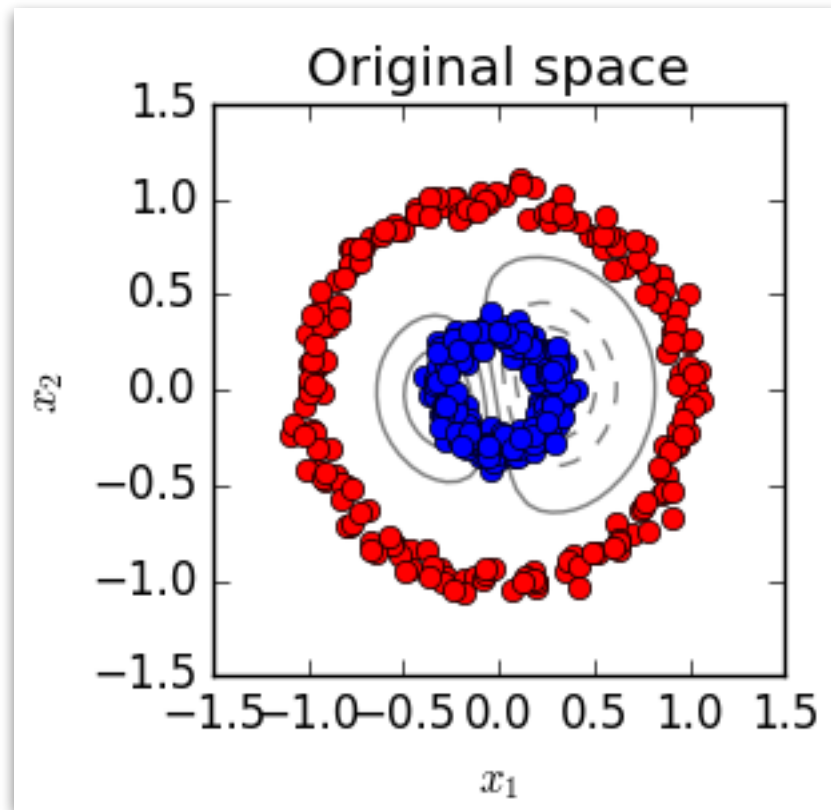
Idea: Use kernels

Linear dimensionality reduction in $\phi(\mathbf{x})$ space



Nonlinear dimensionality reduction in \mathbf{x} space

Kernel PCA



Wrapping up

- PCA is a linear model for dimensionality reduction which finds a mapping to a lower dimensional space that maximizes variance

Wrapping up

- PCA is a linear model for dimensionality reduction which finds a mapping to a lower dimensional space that maximizes variance
- We saw that this is equivalent to performing an eigendecomposition on the covariance matrix of \mathbf{X}

Wrapping up

- PCA is a linear model for dimensionality reduction which finds a mapping to a lower dimensional space that maximizes variance
- We saw that this is equivalent to performing an eigendecomposition on the covariance matrix of \mathbf{X}
- ***Next time*** Auto-encoders and neural compression for non-linear projections