

# Dimensionality Reduction

## Via Principal Component Analysis

Probability and Statistics for Data Science

Carlos Fernandez-Granda



These slides are based on the book [Probability and Statistics for Data Science](#) by Carlos Fernandez-Granda, available for purchase [here](#). A free preprint, videos, code, slides and solutions to exercises are available at <https://www.ps4ds.net>

# Motivation

Data with a large number of features can be difficult to analyze/process

**Solution:** Reduce dimensionality while preserving as much information as possible

Important **preprocessing** step in many applications

# Linear dimensionality reduction

We model data as samples from  $d$ -dimensional random vector  $\tilde{x}$

$\tilde{x}$  has zero mean (otherwise we center by subtracting mean)

For any orthonormal basis  $b_1, \dots, b_d$

$$\tilde{x} = \sum_{i=1}^d \tilde{a}[i] b_i \quad \tilde{a}[i] := b_i^T \tilde{x}$$

$\tilde{a}[1], \tilde{a}[2], \dots, \tilde{a}[d]$  is an equivalent representation of  $\tilde{x}$

Idea: Use only first  $k < d$  coefficients

$$\underset{b_1, \dots, b_k}{\text{approx}}(\tilde{x}) := \sum_{i=1}^k \tilde{a}[i] b_i \quad \tilde{a}[i] := b_i^T \tilde{x}$$

How do we choose  $b_1, \dots, b_k$ ?

$$\underbrace{\sum_{i=1}^d \tilde{a}[i] b_i}_{\tilde{x}} = \underbrace{\sum_{i=1}^k \tilde{a}[i] b_i}_{\text{approx}(\tilde{x})_{b_1, \dots, b_k}} + \underbrace{\sum_{i=k+1}^d \tilde{a}[i] b_i}_{\text{error}}$$

$$\|\tilde{x}\|_2^2 = \left\| \sum_{i=1}^k \tilde{a}[i] b_i \right\|_2^2 + \|\text{error}\|_2^2$$

$$\|\text{error}\|_2^2 = \|\tilde{x}\|_2^2 - \left\| \sum_{i=1}^k \tilde{a}[i] b_i \right\|_2^2$$

$$= \|\tilde{x}\|_2^2 - \sum_{i=1}^k \tilde{a}[i]^2$$

$$= \|\tilde{x}\|_2^2 - \sum_{i=1}^k \left( b_i^T \tilde{x} \right)^2$$

## Mean $\ell_2$ -norm error

$$\begin{aligned}\mathbb{E} \left[ \|\text{error}\|_2^2 \right] &= \mathbb{E} \left[ \|\tilde{x}\|_2^2 \right] - \sum_{i=1}^k \mathbb{E} \left[ \left( b_i^T \tilde{x} \right)^2 \right] \\ &= \mathbb{E} \left[ \|\tilde{x}\|_2^2 \right] - \sum_{i=1}^k \text{Var} \left[ b_i^T \tilde{x} \right]\end{aligned}$$

What  $b_1, \dots, b_k$  maximize directional variance?

# Principal directions

Let  $u_1, \dots, u_d$  be the eigenvectors of covariance matrix  $\Sigma_{\tilde{x}}$

$$u_1 = \arg \max_{\|a\|_2=1} \text{Var}[a^T \tilde{x}]$$

$$u_k = \arg \max_{\|a\|_2=1, a \perp u_1, \dots, u_{k-1}} \text{Var}[a^T \tilde{x}] \quad 2 \leq k \leq d$$

# Principal component analysis

Let  $u_1, \dots, u_d$  be the eigenvectors of covariance matrix  $\Sigma_{\tilde{x}}$

$$\{u_1, \dots, u_k\} = \arg \min_{\substack{\{b_1, \dots, b_k\} \\ \|b_i\|_2=1, 1 \leq i \leq k \\ b_i \perp b_j, i \neq j}} \mathbb{E} \left[ \left\| \tilde{x} - \underset{b_1, \dots, b_k}{\text{approx}}(\tilde{x}) \right\|_2^2 \right]$$

The optimal linear  $k$ -dimensional approximation is

$$\underset{u_1, \dots, u_k}{\text{approx}}(\tilde{x}) = \sum_{i=1}^k \tilde{w}_i u_i \quad \tilde{w}_i := u_i^T \tilde{x}$$



# Wheat seeds

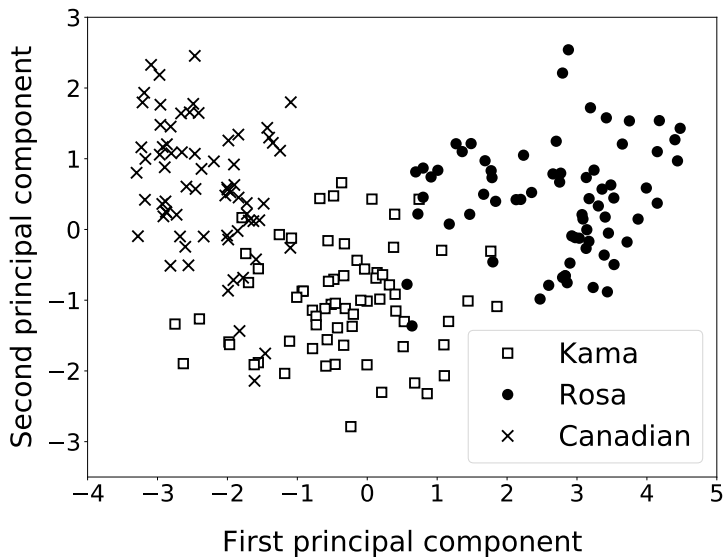
3 varieties: Kama, Rosa and Canadian

Features:

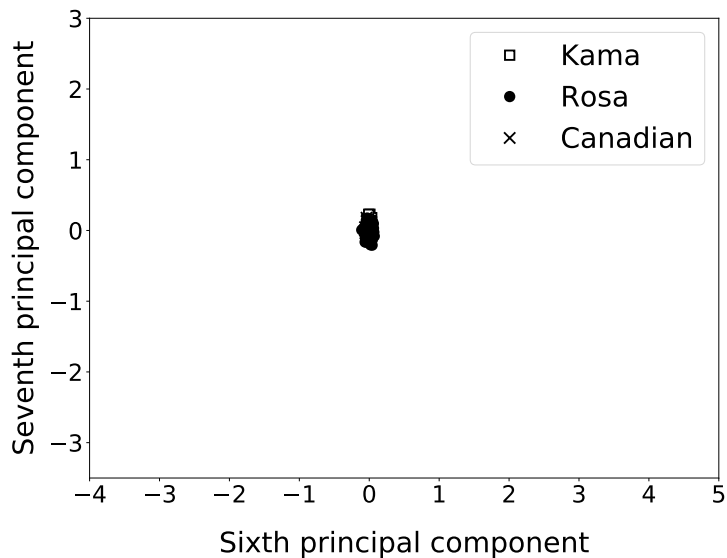
- ▶ Area
- ▶ Perimeter
- ▶ Compactness
- ▶ Length of kernel
- ▶ Width of kernel
- ▶ Asymmetry coefficient
- ▶ Length of kernel groove

**Challenge:** How to visualize the data in two dimensions?

## Two first principal components



## Two last principal components



# Faces

$64 \times 64$  images from 40 subjects

Vectorized images interpreted as vectors in  $\mathbb{R}^{4096}$



Sample mean

## Principal directions

$u_1$



18.8

$u_2$



11.1

$u_3$



6.30

$u_4$



3.95

$u_5$



2.86

## Principal directions

$u_{10}$



1.32

$u_{20}$



0.591

$u_{30}$



0.349

$u_{40}$



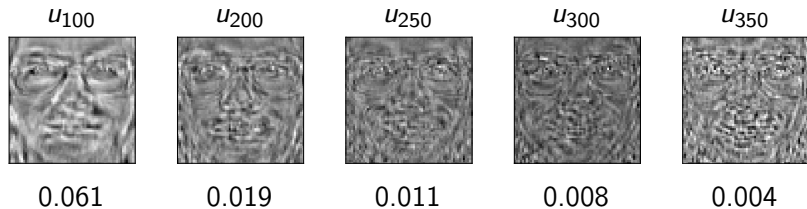
0.217

$u_{50}$



0.162

## Principal directions



# Faces

$$\underset{u_1, \dots, u_5}{\text{approx}}(x_i) := m(X) + \sum_{j=1}^5 w_j[i] u_j \qquad w_j[i] := u_j^T \text{ct}(x_i)$$



$k = 5$

$\text{approx}(x_i)$   
 $u_1, \dots, u_5$

$m(X)$

$w_1$

$u_1$

$w_2$

$u_2$



=



- 1.89



+ 0.92



- 1.08



- 1.51



- 0.73



$w_3$

$u_3$

$w_4$

$u_4$

$w_5$

$u_5$

# Approximation

Original



$k = 5$



$k = 10$



$k = 20$



$k = 30$



$k = 50$



$k = 100$



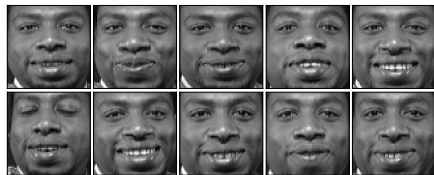
$k = 300$



# Face recognition

Goal: Identify person

Training set:  $\{x_1, y_1\}, \dots, \{x_n, y_n\}$



## Nearest-neighbor classification

$$i^* := \arg \min_{1 \leq i \leq n} \|x_{\text{test}} - x_i\|_2$$

**Cost:**  $\mathcal{O}(nd)$  to classify new point

## Classification in reduced space

Compute sample mean and  $k$  first principal directions  $u_1, \dots, u_k$  from training data

For each test data point  $x_{\text{test}}$

1. Center using training sample mean to obtain  $\text{ct}(x_{\text{test}})$
2. Compute  $k$  principal components

$$w_{\text{test}}[i] := u_i^T \text{ct}(x_{\text{test}})$$

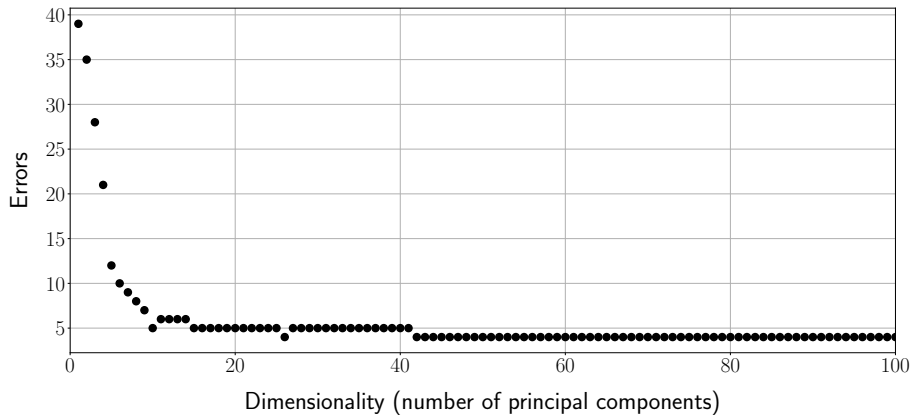
3. Compare to principal components of training data

$$i_{[k]}^* := \arg \min_{1 \leq i \leq n} \|w_{\text{test}} - w_{[1:k]}[i]\|_2$$

Cost reduced to  $\mathcal{O}(nk)$

Computing eigendecomposition is costly, but is done only once

# Performance



$k := 42$

$x_{\text{test}}$



$\text{approx}(x_{\text{test}})$   
 $u_1, \dots, u_k$



$\text{approx}(x_{i[k]}^*)$   
 $u_1, \dots, u_k$



$x_{i[k]}^*$



$k := 42$

$x_{\text{test}}$



$\text{approx}(x_{\text{test}})$   
 $u_1, \dots, u_k$



$\text{approx}(x_{i[k]}^*)$   
 $u_1, \dots, u_k$



$x_{i[k]}^*$



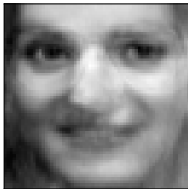


$k := 42$

$x_{\text{test}}$



$\text{approx}(x_{\text{test}})$   
 $u_1, \dots, u_k$



$\text{approx}(x_{i[k]}^*)$   
 $u_1, \dots, u_k$



$x_{i[k]}^*$



$k := 42$

$x_{\text{test}}$



$\text{approx}(x_{\text{test}})$   
 $u_1, \dots, u_k$



$\text{approx}(x_{i[k]}^*)$   
 $u_1, \dots, u_k$



$x_{i[k]}^*$



# Optimality

Let  $u_1, \dots, u_d$  be the eigenvectors of covariance matrix  $\Sigma_{\tilde{x}}$

$$\{u_1, \dots, u_k\} = \arg \min_{\substack{\{b_1, \dots, b_k\} \\ \|b_i\|_2=1, 1 \leq i \leq k \\ b_i \perp b_j, i \neq j}} \mathbb{E} \left[ \left\| \tilde{x} - \underset{b_1, \dots, b_k}{\text{approx}}(\tilde{x}) \right\|_2^2 \right]$$

The **optimal** linear  $k$ -dimensional approximation is

$$\underset{u_1, \dots, u_k}{\text{approx}}(\tilde{x}) := \sum_{i=1}^k \tilde{w}_i u_i \quad \tilde{w}_i := u_i^T \tilde{x}$$

# Proof of optimality

How do we prove this?

$$\{u_1, \dots, u_k\} = \arg \min_{\substack{\{b_1, \dots, b_k\} \\ \|b_i\|_2=1, 1 \leq i \leq k \\ b_i \perp b_j, i \neq j}} \mathbb{E} \left[ \left\| \tilde{x} - \underset{b_1, \dots, b_k}{\text{approx}}(\tilde{x}) \right\|_2^2 \right]$$

## Sum of directional variances

$$\mathbb{E} \left[ \left\| \tilde{x} - \underset{b_1, \dots, b_k}{\text{approx}}(\tilde{x}) \right\|_2^2 \right] = \mathbb{E} \left[ \|\tilde{x}\|_2^2 \right] - \sum_{i=1}^k \text{Var} \left[ b_i^T \tilde{x} \right]$$

We prove that principal directions are optimal by induction on  $k$

$$k = 1$$

By the spectral theorem

$$u_1 = \arg \max_{\|b\|_2=1} \text{Var}[b^T \tilde{x}]$$

## Induction step

We need to show that

$$\{u_1, \dots, u_{k-1}\} = \arg \max_{\substack{\{b_1, \dots, b_{k-1}\} \\ \|b_i\|_2=1, 1 \leq i \leq k-1 \\ b_i \perp b_j, i \neq j}} \sum_{i=1}^{k-1} \text{Var} \left[ b_i^T \tilde{x} \right]$$

implies

$$\{u_1, \dots, u_k\} = \arg \max_{\substack{\{b_1, \dots, b_k\} \\ \|b_i\|_2=1, 1 \leq i \leq k \\ b_i \perp b_j, i \neq j}} \sum_{i=1}^k \text{Var} \left[ b_i^T \tilde{x} \right]$$

## Sum of variances

Fix arbitrary set of  $k$  orthonormal vectors  $b_1, \dots, b_k$

$$\begin{aligned}\sum_{i=1}^k \text{Var} \left[ b_i^T \tilde{x} \right] &= \sum_{i=1}^k \text{E} \left[ \left( b_i^T \tilde{x} \right)^2 \right] \\&= \text{E} \left[ \sum_{i=1}^k \left( b_i^T \tilde{x} \right)^2 \right] \\&= \text{E} \left[ \left\| \sum_{i=1}^k b_i^T \tilde{x} b_i \right\|_2^2 \right] \\&= \text{E} \left[ \left\| \mathcal{P}_{\mathcal{S}} \tilde{x} \right\|_2^2 \right] \quad \mathcal{S} := \text{span}(b_1, \dots, b_k)\end{aligned}$$



## Key insight

For any orthonormal basis  $a_1, \dots, a_k$  of  $\mathcal{S}$

$$\begin{aligned}\mathcal{P}_{\mathcal{S}} \tilde{x} &:= \sum_{i=1}^k b_i^T \tilde{x} b_i \\ &= \sum_{i=1}^k a_i^T \tilde{x} a_i\end{aligned}$$

We need to choose wisely!

## Orthogonal vector

$\mathcal{S}$  has dimension  $k$

There is at least one vector  $a_{\perp} \in \mathcal{S}$  orthogonal to  $u_1, \dots, u_{k-1}$

Can we have

$$\text{Var}[u_k^T \tilde{x}] < \text{Var}[a_{\perp}^T \tilde{x}] \quad ?$$

No!

$$u_k = \arg \max_{\|a\|_2=1, a \perp u_1, \dots, u_{k-1}} \text{Var}[a^T \tilde{x}]$$

## Wise choice

Set  $a_k := a_\perp$

Choose  $a_1, \dots, a_{k-1}$  so  $a_1, \dots, a_k$  span  $\mathcal{S}$

$$\text{Var}[u_k^T \tilde{x}] \geq \text{Var}[a_\perp^T \tilde{x}] = \text{Var}[a_k^T \tilde{x}]$$

$$\sum_{i=1}^{k-1} \text{Var}[u_i^T \tilde{x}] \geq \sum_{i=1}^{k-1} \text{Var}[a_i^T \tilde{x}] \quad \text{by induction hypothesis}$$

$$\begin{aligned} \sum_{i=1}^k \text{Var}[u_i^T \tilde{x}] &\geq \sum_{i=1}^k \text{Var}[a_i^T \tilde{x}] \\ &= \mathbb{E} \left[ \|\mathcal{P}_{\mathcal{S}} \tilde{x}\|_2^2 \right] = \sum_{i=1}^k \text{Var}[b_i^T \tilde{x}] \end{aligned}$$

# What have we learned?

Definition of dimensionality reduction

How to perform linear dimensionality reduction via PCA

Proof of optimality for mean  $\ell_2$ -norm error