

# Introduction to Machine Learning

## Feature Generation

Andres Mendez-Vazquez

June 14, 2020

# Outline

## 1 Fisher Linear Discriminant

- Introduction
- The Rotation Idea
- Solution
  - Scatter measure
- The Cost Function

## 2 Principal Components and Singular Value Decomposition

- Introduction
- Principal Component Analysis AKA Karhunen-Loeve Transform
  - Projecting the Data
  - Lagrange Multipliers
  - The Process
  - Example
- Singular Value Decomposition
  - Introduction
  - Image Compression

# Outline

## 1 Fisher Linear Discriminant

- Introduction
  - The Rotation Idea
  - Solution
    - Scatter measure
  - The Cost Function

## 2 Principal Components and Singular Value Decomposition

- Introduction
- Principal Component Analysis AKA Karhunen-Loeve Transform
  - Projecting the Data
  - Lagrange Multipliers
  - The Process
  - Example
- Singular Value Decomposition
  - Introduction
  - Image Compression

# What do we want?

## What

- Given a set of measurements, the goal is to discover compact and informative representations of the obtained data.

## Our Approach

- We want to “squeeze” in a relatively small number of features, leading to a reduction of the necessary feature space dimension.

## Properties

- Thus removing information redundancies - Usually produced and the measurement.

# What do we want?

## What

- Given a set of measurements, the goal is to discover compact and informative representations of the obtained data.

## Our Approach

- We want to “squeeze” in a relatively small number of features, leading to a reduction of the necessary feature space dimension.

## Properties

- Thus removing information redundancies - Usually produced and the measurement.

# What do we want?

## What

- Given a set of measurements, the goal is to discover compact and informative representations of the obtained data.

## Our Approach

- We want to “squeeze” in a relatively small number of features, leading to a reduction of the necessary feature space dimension.

## Properties

- Thus removing information redundancies - Usually produced and the measurement.

# What Methods we will see?

## Fisher Linear Discriminant

- 1 Squeezing to the maximum.
- 2 From Many to One Dimension

## Principal Component Analysis

- Not so much squeezing
- You are willing to lose some information

# What Methods we will see?

## Fisher Linear Discriminant

- 1 Squeezing to the maximum.
- 2 From Many to One Dimension

## Principal Component Analysis

- 1 Not so much squeezing
- 2 You are willing to lose some information



# Outline

## 1 Fisher Linear Discriminant

- Introduction
- **The Rotation Idea**
- Solution
  - Scatter measure
- The Cost Function

## 2 Principal Components and Singular Value Decomposition

- Introduction
- Principal Component Analysis AKA Karhunen-Loeve Transform
  - Projecting the Data
  - Lagrange Multipliers
  - The Process
  - Example
- Singular Value Decomposition
  - Introduction
  - Image Compression

# Rotation

## Projecting

Projecting well-separated samples onto an arbitrary line usually produces a confused mixture of samples from all of the classes and thus produces poor recognition performance.

## Something Notable

However, moving and rotating the line around might result in an orientation for which the projected samples are well separated.

## Best Linear Discriminant (LSD)

It is a discriminant analysis seeking directions that are efficient for discriminating binary classification problem.

# Rotation

## Projecting

Projecting well-separated samples onto an arbitrary line usually produces a confused mixture of samples from all of the classes and thus produces poor recognition performance.

## Something Notable

However, moving and rotating the line around might result in an orientation for which the projected samples are well separated.

## Linear discriminant analysis (LDA)

It is a discriminant analysis seeking directions that are efficient for discriminating binary classification problem.

# Rotation

## Projecting

Projecting well-separated samples onto an arbitrary line usually produces a confused mixture of samples from all of the classes and thus produces poor recognition performance.

## Something Notable

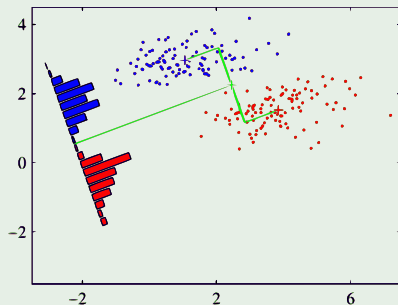
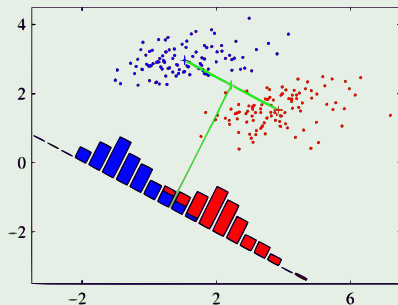
However, moving and rotating the line around might result in an orientation for which the projected samples are well separated.

## Fisher linear discriminant (FLD)

It is a discriminant analysis seeking directions that are efficient for discriminating binary classification problem.

# Example

## Example - From Left to Right the Improvement



This is actually coming from...

## Classifier as

A machine for dimensionality reduction.

### Initial Setup

We have:

- $N$   $d$ -dimensional samples  $x_1, x_2, \dots, x_N$
- $N_i$  is the number of samples in class  $C_i$  for  $i=1,2$ .

Then, we ask for the projection of each  $x_i$  into the line by means of

$$y_i = w^T x_i \quad (1)$$

This is actually coming from...

## Classifier as

A machine for dimensionality reduction.

## Initial Setup

We have:

- $N$   $d$ -dimensional samples  $x_1, x_2, \dots, x_N$
- $N_i$  is the number of samples in class  $C_i$  for  $i=1,2$ .

Then, we ask for the projection of each  $x_i$  into the line by means of

$$y_i = w^T x_i \quad (1)$$

This is actually coming from...

## Classifier as

A machine for dimensionality reduction.

## Initial Setup

We have:

- $N$   $d$ -dimensional samples  $x_1, x_2, \dots, x_N$
- $N_i$  is the number of samples in class  $C_i$  for  $i=1,2$ .

Then, we ask for the projection of each  $x_i$  into the line by means of

$$y_i = \mathbf{w}^T \mathbf{x}_i \quad (1)$$



# Outline

## 1 Fisher Linear Discriminant

- Introduction
- The Rotation Idea
- **Solution**
  - Scatter measure
  - The Cost Function

## 2 Principal Components and Singular Value Decomposition

- Introduction
- Principal Component Analysis AKA Karhunen-Loeve Transform
  - Projecting the Data
  - Lagrange Multipliers
  - The Process
  - Example
- Singular Value Decomposition
  - Introduction
  - Image Compression

# Use the mean of each Class

Then

Select  $w$  such that class separation is maximized

We then define the mean sample for each class

$$\bullet C_1 \Rightarrow m_1 = \frac{1}{N_1} \sum_{i=1}^{N_1} x_i$$

$$\bullet C_2 \Rightarrow m_2 = \frac{1}{N_2} \sum_{i=1}^{N_2} x_i$$

OK!! This is giving us a measure of distance

Thus, we want to maximize the distance the projected means:

$$m_1 - m_2 = w^T (m_1 - m_2) \quad (2)$$

where  $m_k = w^T m_k$  for  $k = 1, 2$ .

# Use the mean of each Class

Then

Select  $w$  such that class separation is maximized

We then define the mean sample for each class

$$① \quad C_1 \Rightarrow m_1 = \frac{1}{N_1} \sum_{i=1}^{N_1} x_i$$

$$② \quad C_2 \Rightarrow m_2 = \frac{1}{N_2} \sum_{i=1}^{N_2} x_i$$

OK!! This is giving us a measure of distance

Thus, we want to maximize the distance the projected means:

$$m_1 - m_2 = w^T (m_1 - m_2) \quad (2)$$

where  $m_k = w^T m_k$  for  $k = 1, 2$ .

## Use the mean of each Class

Then

Select  $\mathbf{w}$  such that class separation is maximized

We then define the mean sample for each class

$$\textcircled{1} \quad C_1 \Rightarrow \mathbf{m}_1 = \frac{1}{N_1} \sum_{i=1}^{N_1} \mathbf{x}_i$$

$$\textcircled{2} \quad C_2 \Rightarrow \mathbf{m}_2 = \frac{1}{N_2} \sum_{i=1}^{N_2} \mathbf{x}_i$$

Ok!!! This is giving us a measure of distance

Thus, we want to maximize the distance the projected means:

$$m_1 - m_2 = \mathbf{w}^T (\mathbf{m}_1 - \mathbf{m}_2) \quad (2)$$

where  $m_k = \mathbf{w}^T \mathbf{m}_k$  for  $k = 1, 2$ .

However

We could simply seek

$$\begin{aligned} \max \mathbf{w}^T (\mathbf{m}_1 - \mathbf{m}_2) \\ \text{s.t. } \sum_{i=1}^d w_i = 1 \end{aligned}$$

After all

We do not care about the magnitude of  $\mathbf{w}$ .

However

We could simply seek

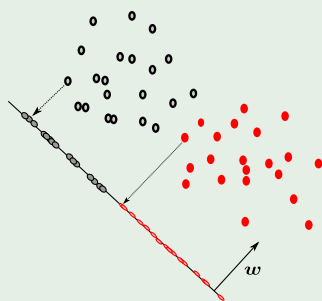
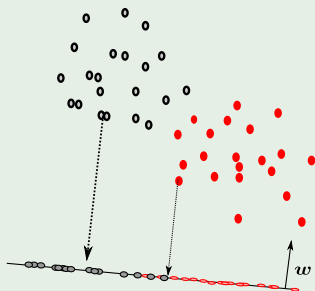
$$\begin{aligned} \max \mathbf{w}^T (\mathbf{m}_1 - \mathbf{m}_2) \\ \text{s.t. } \sum_{i=1}^d w_i = 1 \end{aligned}$$

After all

We do not care about the magnitude of  $\mathbf{w}$ .

# Example

Here, we have the problem



# Outline

## 1 Fisher Linear Discriminant

- Introduction
- The Rotation Idea
- **Solution**
  - **Scatter measure**
- The Cost Function

## 2 Principal Components and Singular Value Decomposition

- Introduction
- Principal Component Analysis AKA Karhunen-Loeve Transform
  - Projecting the Data
  - Lagrange Multipliers
  - The Process
  - Example
- Singular Value Decomposition
  - Introduction
  - Image Compression



# Fixing the Problem

To obtain good separation of the projected data

The difference between the means should be large relative to some measure of the standard deviations for each class.

We define a SCATTER measure (Based in the Sample Variance)

$$s_k^2 = \sum_{x_i \in C_k} (w^T x_i - m_k)^2 = \sum_{y_i = w^T x_i \in C_k} (y_i - m_k)^2 \quad (3)$$

We define then within-class variance for the whole data

$$s_1^2 + s_2^2 \quad (4)$$

# Fixing the Problem

To obtain good separation of the projected data

The difference between the means should be large relative to some measure of the standard deviations for each class.

We define a SCATTER measure (Based in the Sample Variance)

$$s_k^2 = \sum_{\mathbf{x}_i \in C_k} \left( \mathbf{w}^T \mathbf{x}_i - m_k \right)^2 = \sum_{y_i = \mathbf{w}^T \mathbf{x}_i \in C_k} (y_i - m_k)^2 \quad (3)$$

We define then within-class variance for the whole data

$$s_1^2 + s_2^2 \quad (4)$$

# Fixing the Problem

To obtain good separation of the projected data

The difference between the means should be large relative to some measure of the standard deviations for each class.

We define a SCATTER measure (Based in the Sample Variance)

$$s_k^2 = \sum_{\mathbf{x}_i \in C_k} \left( \mathbf{w}^T \mathbf{x}_i - m_k \right)^2 = \sum_{y_i = \mathbf{w}^T \mathbf{x}_i \in C_k} (y_i - m_k)^2 \quad (3)$$

We define then within-class variance for the whole data

$$s_1^2 + s_2^2 \quad (4)$$

# Outline

## 1 Fisher Linear Discriminant

- Introduction
- The Rotation Idea
- Solution
  - Scatter measure
- **The Cost Function**

## 2 Principal Components and Singular Value Decomposition

- Introduction
- Principal Component Analysis AKA Karhunen-Loeve Transform
  - Projecting the Data
  - Lagrange Multipliers
  - The Process
  - Example
- Singular Value Decomposition
  - Introduction
  - Image Compression

## Finally, a Cost Function

The between-class variance

$$(m_1 - m_2)^2 \quad (5)$$

The Fisher criterion

$$\frac{\text{between-class variance}}{\text{within-class variance}} \quad (6)$$

Finally

$$J(w) = \frac{(m_1 - m_2)^2}{s_1^2 + s_2^2} \quad (7)$$

## Finally, a Cost Function

The between-class variance

$$(m_1 - m_2)^2 \quad (5)$$

The Fisher criterion

$$\frac{\text{between-class variance}}{\text{within-class variance}} \quad (6)$$

Finally

$$J(w) = \frac{(m_1 - m_2)^2}{s_1^2 + s_2^2} \quad (7)$$

## Finally, a Cost Function

The between-class variance

$$(m_1 - m_2)^2 \quad (5)$$

The Fisher criterion

$$\frac{\text{between-class variance}}{\text{within-class variance}} \quad (6)$$

Finally

$$J(\mathbf{w}) = \frac{(m_1 - m_2)^2}{s_1^2 + s_2^2} \quad (7)$$

# We use a transformation to simplify our life

## First

$$J(w) = \frac{(w^T m_1 - w^T m_2)^2}{\sum_{y_i=w^T x_i \in C_1} (y_i - m_k)^2 + \sum_{y_i=w^T x_i \in C_2} (y_i - m_k)^2}$$

## Second

$$= \frac{(w^T m_1 - w^T m_2) (w^T m_1 - w^T m_2)^T}{\sum_{y_i=w^T x_i \in C_1} (w^T x_i - m_k) (w^T x_i - m_k)^T + \sum_{y_i=w^T x_i \in C_2} (w^T x_i - m_k) (w^T x_i - m_k)^T}$$

## Third

$$= \frac{w^T (m_1 - m_2) (w^T (m_1 - m_2))^T}{\sum_{y_i=w^T x_i \in C_1} w^T (x_i - m_1) (w^T (x_i - m_1))^T + \sum_{y_i=w^T x_i \in C_2} w^T (x_i - m_2) (w^T (x_i - m_2))^T}$$



# We use a transformation to simplify our life

## First

$$J(\mathbf{w}) = \frac{(\mathbf{w}^T \mathbf{m}_1 - \mathbf{w}^T \mathbf{m}_2)^2}{\sum_{y_i = \mathbf{w}^T \mathbf{x}_i \in C_1} (y_i - m_k)^2 + \sum_{y_i = \mathbf{w}^T \mathbf{x}_i \in C_2} (y_i - m_k)^2}$$

## Second

$$= \frac{(\mathbf{w}^T \mathbf{m}_1 - \mathbf{w}^T \mathbf{m}_2) (\mathbf{w}^T \mathbf{m}_1 - \mathbf{w}^T \mathbf{m}_2)^T}{\sum_{y_i = \mathbf{w}^T \mathbf{x}_i \in C_1} (\mathbf{w}^T \mathbf{x}_i - m_k) (\mathbf{w}^T \mathbf{x}_i - m_k)^T + \sum_{y_i = \mathbf{w}^T \mathbf{x}_i \in C_2} (\mathbf{w}^T \mathbf{x}_i - m_k) (\mathbf{w}^T \mathbf{x}_i - m_k)^T}$$

## Third

$$= \frac{\mathbf{w}^T (\mathbf{m}_1 - \mathbf{m}_2) (\mathbf{w}^T (\mathbf{m}_1 - \mathbf{m}_2))^T}{\sum_{y_i = \mathbf{w}^T \mathbf{x}_i \in C_1} \mathbf{w}^T (\mathbf{x}_i - m_1) (\mathbf{w}^T (\mathbf{x}_i - m_1))^T + \sum_{y_i = \mathbf{w}^T \mathbf{x}_i \in C_2} \mathbf{w}^T (\mathbf{x}_i - m_2) (\mathbf{w}^T (\mathbf{x}_i - m_2))^T}$$

# We use a transformation to simplify our life

## First

$$J(w) = \frac{(w^T m_1 - w^T m_2)^2}{\sum_{y_i = w^T x_i \in C_1} (y_i - m_k)^2 + \sum_{y_i = w^T x_i \in C_2} (y_i - m_k)^2}$$

## Second

$$= \frac{(w^T m_1 - w^T m_2) (w^T m_1 - w^T m_2)^T}{\sum_{y_i = w^T x_i \in C_1} (w^T x_i - m_k) (w^T x_i - m_k)^T + \sum_{y_i = w^T x_i \in C_2} (w^T x_i - m_k) (w^T x_i - m_k)^T}$$

## Third

$$= \frac{w^T (m_1 - m_2) (w^T (m_1 - m_2))^T}{\sum_{y_i = w^T x_i \in C_1} w^T (x_i - m_1) (w^T (x_i - m_1))^T + \sum_{y_i = w^T x_i \in C_2} w^T (x_i - m_2) (w^T (x_i - m_2))^T}$$

# Transformation

## Fourth

$$= \frac{w^T (m_1 - m_2) (m_1 - m_2)^T w}{\sum_{y_i = w^T x_i \in C_1} w^T (x_i - m_1) (x_i - m_1)^T w + \sum_{y_i = w^T x_i \in C_2} w^T (x_i - m_2) (x_i - m_2)^T w}$$

## Fifth

$$= \frac{w^T (m_1 - m_2) (m_1 - m_2)^T w}{w^T \left[ \sum_{y_i = w^T x_i \in C_1} (x_i - m_1) (x_i - m_1)^T + \sum_{y_i = w^T x_i \in C_2} (x_i - m_2) (x_i - m_2)^T \right] w}$$

## Now Rename

$$J(w) = \frac{w^T S_B w}{w^T S_w w} \quad (8)$$

# Transformation

## Fourth

$$= \frac{w^T (m_1 - m_2) (m_1 - m_2)^T w}{\sum_{y_i=w^T x_i \in C_1} w^T (x_i - m_1) (x_i - m_1)^T w + \sum_{y_i=w^T x_i \in C_2} w^T (x_i - m_2) (x_i - m_2)^T w}$$

## Fifth

$$= \frac{w^T (m_1 - m_2) (m_1 - m_2)^T w}{w^T \left[ \sum_{y_i=w^T x_i \in C_1} (x_i - m_1) (x_i - m_1)^T + \sum_{y_i=w^T x_i \in C_2} (x_i - m_2) (x_i - m_2)^T \right] w}$$

## Now Define

$$J(w) = \frac{w^T S_B w}{w^T S_w w} \quad (8)$$

# Transformation

## Fourth

$$= \frac{w^T (m_1 - m_2) (m_1 - m_2)^T w}{\sum_{y_i = w^T x_i \in C_1} w^T (x_i - m_1) (x_i - m_1)^T w + \sum_{y_i = w^T x_i \in C_2} w^T (x_i - m_2) (x_i - m_2)^T w}$$

## Fifth

$$= \frac{w^T (m_1 - m_2) (m_1 - m_2)^T w}{w^T \left[ \sum_{y_i = w^T x_i \in C_1} (x_i - m_1) (x_i - m_1)^T + \sum_{y_i = w^T x_i \in C_2} (x_i - m_2) (x_i - m_2)^T \right] w}$$

## Now Rename

$$J(w) = \frac{w^T S_B w}{w^T S_w w} \quad (8)$$

## Derive with respect to $w$

Thus

$$\frac{dJ(w)}{dw} = \frac{d(w^T S_B w) (w^T S_w w)^{-1}}{dw} = 0 \quad (9)$$

Then

$$\frac{dJ(w)}{dw} = (S_B w + S_B^T w) (w^T S_w w)^{-1} - (w^T S_B w) (w^T S_w w)^{-2} (S_w w + S_w^T w) = 0 \quad (10)$$

Now, because the symmetry in  $S_B$  and  $S_w$ ,

$$\frac{dJ(w)}{dw} = \frac{S_B}{(w^T S_w w)} - \frac{w^T S_B w S_w w}{(w^T S_w w)^2} = 0 \quad (11)$$

Derive with respect to  $w$

Thus

$$\frac{dJ(w)}{dw} = \frac{d(w^T S_B w) (w^T S_w w)^{-1}}{dw} = 0 \quad (9)$$

Then

$$\frac{dJ(w)}{dw} = (S_B w + S_B^T w) (w^T S_w w)^{-1} - (w^T S_B w) (w^T S_w w)^{-2} (S_w w + S_w^T w) = 0 \quad (10)$$

Now, because the symmetry in  $S_w$  and  $S_B$ ,

$$\frac{dJ(w)}{dw} = \frac{S_B}{(w^T S_w w)} - \frac{w^T S_B w S_w w}{(w^T S_w w)^2} = 0 \quad (11)$$

Derive with respect to  $w$

Thus

$$\frac{dJ(w)}{dw} = \frac{d(w^T S_B w) (w^T S_w w)^{-1}}{dw} = 0 \quad (9)$$

Then

$$\frac{dJ(w)}{dw} = (S_B w + S_B^T w) (w^T S_w w)^{-1} - (w^T S_B w) (w^T S_w w)^{-2} (S_w w + S_w^T w) = 0 \quad (10)$$

Now because the symmetry in  $S_B$  and  $S_w$

$$\frac{dJ(w)}{dw} = \frac{S_B}{(w^T S_w w)} - \frac{w^T S_B w S_w w}{(w^T S_w w)^2} = 0 \quad (11)$$



Derive with respect to  $w$

Thus

$$\frac{dJ(w)}{dw} = \frac{S_B}{(w^T S_w w)} - \frac{w^T S_B w S_w w}{(w^T S_w w)^2} = 0 \quad (12)$$

Then

$$(w^T S_w w) S_B w = (w^T S_B w) S_w w \quad (13)$$

Derive with respect to  $w$

Thus

$$\frac{dJ(w)}{dw} = \frac{S_B}{(w^T S_w w)} - \frac{w^T S_B w S_w w}{(w^T S_w w)^2} = 0 \quad (12)$$

Then

$$(w^T S_w w) S_B w = (w^T S_B w) S_w w \quad (13)$$

## Now, Several Tricks!!!

### First

$$\mathbf{S}_B \mathbf{w} = (\mathbf{m}_1 - \mathbf{m}_2) (\mathbf{m}_1 - \mathbf{m}_2)^T \mathbf{w} = \alpha (\mathbf{m}_1 - \mathbf{m}_2) \quad (14)$$

Where  $\alpha = (\mathbf{m}_1 - \mathbf{m}_2)^T \mathbf{w}$  is a simple constant

It means that  $\mathbf{S}_B \mathbf{w}$  is always in the direction  $\mathbf{m}_1 - \mathbf{m}_2$ !!!

In addition

$\mathbf{w}^T \mathbf{S}_w \mathbf{w}$  and  $\mathbf{w}^T \mathbf{S}_B \mathbf{w}$  are constants

## Now, Several Tricks!!!

First

$$\mathbf{S}_B \mathbf{w} = (\mathbf{m}_1 - \mathbf{m}_2) (\mathbf{m}_1 - \mathbf{m}_2)^T \mathbf{w} = \alpha (\mathbf{m}_1 - \mathbf{m}_2) \quad (14)$$

Where  $\alpha = (\mathbf{m}_1 - \mathbf{m}_2)^T \mathbf{w}$  is a simple constant

It means that  $\mathbf{S}_B \mathbf{w}$  is always in the direction  $\mathbf{m}_1 - \mathbf{m}_2$ !!!

In addition

$\mathbf{w}^T \mathbf{S}_W \mathbf{w}$  and  $\mathbf{w}^T \mathbf{S}_B \mathbf{w}$  are constants

## Now, Several Tricks!!!

First

$$\mathbf{S}_B \mathbf{w} = (\mathbf{m}_1 - \mathbf{m}_2) (\mathbf{m}_1 - \mathbf{m}_2)^T \mathbf{w} = \alpha (\mathbf{m}_1 - \mathbf{m}_2) \quad (14)$$

Where  $\alpha = (\mathbf{m}_1 - \mathbf{m}_2)^T \mathbf{w}$  is a simple constant

It means that  $\mathbf{S}_B \mathbf{w}$  is always in the direction  $\mathbf{m}_1 - \mathbf{m}_2$ !!!

In addition

$\mathbf{w}^T \mathbf{S}_w \mathbf{w}$  and  $\mathbf{w}^T \mathbf{S}_B \mathbf{w}$  are constants

## Now, Several Tricks!!!

Finally

$$\mathbf{S}_w \mathbf{w} \propto (\mathbf{m}_1 - \mathbf{m}_2) \Rightarrow \mathbf{w} \propto \mathbf{S}_w^{-1} (\mathbf{m}_1 - \mathbf{m}_2) \quad (15)$$

## Now, Several Tricks!!!

### Finally

$$\mathbf{S}_w \mathbf{w} \propto (\mathbf{m}_1 - \mathbf{m}_2) \Rightarrow \mathbf{w} \propto \mathbf{S}_w^{-1} (\mathbf{m}_1 - \mathbf{m}_2) \quad (15)$$

### Once the data is transformed into $y_i$

- Use a threshold  $y_0 \Rightarrow x \in C_1$  iff  $y(x) \geq y_0$  or  $x \in C_2$  iff  $y(x) < y_0$
- Or ML with a Gaussian can be used to classify the new transformed data using a Naive Bayes (Central Limit Theorem and  $y = \mathbf{w}^T \mathbf{x}$  sum of random variables).

## Now, Several Tricks!!!

### Finally

$$\mathbf{S}_w \mathbf{w} \propto (\mathbf{m}_1 - \mathbf{m}_2) \Rightarrow \mathbf{w} \propto \mathbf{S}_w^{-1} (\mathbf{m}_1 - \mathbf{m}_2) \quad (15)$$

### Once the data is transformed into $y_i$

- Use a threshold  $y_0 \Rightarrow x \in C_1$  iff  $y(x) \geq y_0$  or  $x \in C_2$  iff  $y(x) < y_0$
- Or ML with a Gaussian can be used to classify the new transformed data using a Naive Bayes (Central Limit Theorem and  $y = \mathbf{w}^T \mathbf{x}$  sum of random variables).



Please

Your Reading Material, it is about the Multiclass

4.1.6 Fisher's discriminant for multiple classes AT "Pattern Recognition"  
by Bishop

# Outline

## 1 Fisher Linear Discriminant

- Introduction
- The Rotation Idea
- Solution
  - Scatter measure
- The Cost Function

## 2 Principal Components and Singular Value Decomposition

- **Introduction**
- Principal Component Analysis AKA Karhunen-Loeve Transform
  - Projecting the Data
  - Lagrange Multipliers
  - The Process
  - Example
- Singular Value Decomposition
  - Introduction
  - Image Compression

# Did you noticed?

That Rotations really do not exist

- Actually, they are mappings or projections in linear algebra

Thus, Can we get more powerful mappings?

- To obtain better features

Clearly, Yes

- For example, Principal Components or Singular Value Decomposition's

# Did you noticed?

That Rotations really do not exist

- Actually, they are mappings or projections in linear algebra

Thus, Can we get more powerful mappings?

- To obtain better features

Clearly, Yes

- For example, Principal Components or Singular Value Decomposition's

# Did you noticed?

That Rotations really do not exist

- Actually, they are mappings or projections in linear algebra

Thus, Can we get more powerful mappings?

- To obtain better features

Clearly... Yes

- For example, Principal Components or Singular Value Decomposition's

# Outline

## 1 Fisher Linear Discriminant

- Introduction
- The Rotation Idea
- Solution
  - Scatter measure
- The Cost Function

## 2 Principal Components and Singular Value Decomposition

- Introduction
- **Principal Component Analysis AKA Karhunen-Loeve Transform**
  - Projecting the Data
  - Lagrange Multipliers
  - The Process
  - Example
- Singular Value Decomposition
  - Introduction
  - Image Compression

# Also Known as Karhunen-Loeve Transform

## Setup

- Consider a data set of observations  $\{x_n\}$  with  $n = 1, 2, \dots, N$  and  $x_n \in R^d$ .

## Goal

Project data onto space with dimensionality  $m < d$  (We assume  $m$  is given)

# Also Known as Karhunen-Loeve Transform

## Setup

- Consider a data set of observations  $\{x_n\}$  with  $n = 1, 2, \dots, N$  and  $x_n \in R^d$ .

## Goal

Project data onto space with dimensionality  $m < d$  (We assume  $m$  is given)



# Dimensional Variance

Remember the Variance Sample in  $\mathbb{R}$

$$VAR(X) = \frac{\sum_{i=1}^N (x_i - \bar{x})^2}{N - 1} \quad (16)$$

You can do the same in the case of two variables  $X$  and  $Y$

$$COV(x, y) = \frac{\sum_{i=1}^N (x_i - \bar{x})(y_i - \bar{y})}{N - 1} \quad (17)$$

# Dimensional Variance

Remember the Variance Sample in  $\mathbb{R}$

$$VAR(X) = \frac{\sum_{i=1}^N (x_i - \bar{x})^2}{N - 1} \quad (16)$$

You can do the same in the case of two variables  $X$  and  $Y$

$$COV(x, y) = \frac{\sum_{i=1}^N (x_i - \bar{x})(y_i - \bar{y})}{N - 1} \quad (17)$$

## Now, Define

Given the data

$$\mathbf{x}_1, \mathbf{x}_2, \dots, \mathbf{x}_N \quad (18)$$

where  $\mathbf{x}_i$  is a column vector

Construct the sample mean

$$\bar{\mathbf{x}} = \frac{1}{N} \sum_{i=1}^N \mathbf{x}_i \quad (19)$$

Center data

$$\mathbf{x}_1 - \bar{\mathbf{x}}, \mathbf{x}_2 - \bar{\mathbf{x}}, \dots, \mathbf{x}_N - \bar{\mathbf{x}} \quad (20)$$

## Now, Define

Given the data

$$\mathbf{x}_1, \mathbf{x}_2, \dots, \mathbf{x}_N \quad (18)$$

where  $\mathbf{x}_i$  is a column vector

Construct the sample mean

$$\bar{\mathbf{x}} = \frac{1}{N} \sum_{i=1}^N \mathbf{x}_i \quad (19)$$

Center data

$$\mathbf{x}_1 - \bar{\mathbf{x}}, \mathbf{x}_2 - \bar{\mathbf{x}}, \dots, \mathbf{x}_N - \bar{\mathbf{x}} \quad (20)$$

## Now, Define

Given the data

$$\mathbf{x}_1, \mathbf{x}_2, \dots, \mathbf{x}_N \quad (18)$$

where  $\mathbf{x}_i$  is a column vector

Construct the sample mean

$$\bar{\mathbf{x}} = \frac{1}{N} \sum_{i=1}^N \mathbf{x}_i \quad (19)$$

Center data

$$\mathbf{x}_1 - \bar{\mathbf{x}}, \mathbf{x}_2 - \bar{\mathbf{x}}, \dots, \mathbf{x}_N - \bar{\mathbf{x}} \quad (20)$$

# Build the Sample Mean

## The Covariance Matrix

$$S = \frac{1}{N-1} \sum_{i=1}^N (\mathbf{x}_i - \bar{\mathbf{x}}) (\mathbf{x}_i - \bar{\mathbf{x}})^T \quad (21)$$

### Properties

- The  $ij$ th value of  $S$  is equivalent to  $\sigma_{ij}^2$ .
- The  $ii$ th value of  $S$  is equivalent to  $\sigma_{ii}^2$ .

# Build the Sample Mean

## The Covariance Matrix

$$S = \frac{1}{N-1} \sum_{i=1}^N (\mathbf{x}_i - \bar{\mathbf{x}}) (\mathbf{x}_i - \bar{\mathbf{x}})^T \quad (21)$$

## Properties

- 1 The  $ij$ th value of  $S$  is equivalent to  $\sigma_{ij}^2$ .
- 2 The  $ii$ th value of  $S$  is equivalent to  $\sigma_{ii}^2$ .

# Outline

## 1 Fisher Linear Discriminant

- Introduction
- The Rotation Idea
- Solution
  - Scatter measure
- The Cost Function

## 2 Principal Components and Singular Value Decomposition

- Introduction
- Principal Component Analysis AKA Karhunen-Loeve Transform
  - Projecting the Data
    - Lagrange Multipliers
    - The Process
    - Example
- Singular Value Decomposition
  - Introduction
  - Image Compression



## Using $S$ to Project Data

For this we use a  $\mathbf{u}_1$

- with  $\mathbf{u}_1^T \mathbf{u}_1 = 1$ , an orthonormal vector

Question

- What is the Sample Variance of the Projected Data?

# Using $S$ to Project Data

For this we use a  $\mathbf{u}_1$

- with  $\mathbf{u}_1^T \mathbf{u}_1 = 1$ , an orthonormal vector

## Question

- What is the Sample Variance of the Projected Data?

# Outline

## 1 Fisher Linear Discriminant

- Introduction
- The Rotation Idea
- Solution
  - Scatter measure
- The Cost Function

## 2 Principal Components and Singular Value Decomposition

- Introduction
- Principal Component Analysis AKA Karhunen-Loeve Transform
  - Projecting the Data
  - Lagrange Multipliers
  - The Process
  - Example
- Singular Value Decomposition
  - Introduction
  - Image Compression

Thus we have

Variance of the projected data

$$\frac{1}{N-1} \sum_{i=1}^N [\mathbf{u}_1 \mathbf{x}_i - \mathbf{u}_1 \bar{\mathbf{x}}] = \mathbf{u}_1^T S \mathbf{u}_1 \quad (22)$$

Use Lagrange Multipliers to Maximize

$$\mathbf{u}_1^T S \mathbf{u}_1 + \lambda_1 (1 - \mathbf{u}_1^T \mathbf{u}_1) \quad (23)$$

Thus we have

Variance of the projected data

$$\frac{1}{N-1} \sum_{i=1}^N [\mathbf{u}_1 \mathbf{x}_i - \mathbf{u}_1 \bar{\mathbf{x}}] = \mathbf{u}_1^T S \mathbf{u}_1 \quad (22)$$

Use Lagrange Multipliers to Maximize

$$\mathbf{u}_1^T S \mathbf{u}_1 + \lambda_1 (1 - \mathbf{u}_1^T \mathbf{u}_1) \quad (23)$$

Derive by  $\mathbf{u}_1$

We get

$$S\mathbf{u}_1 = \lambda_1\mathbf{u}_1 \quad (24)$$

Then

$\mathbf{u}_1$  is an eigenvector of  $S$ .

If we left-multiply by  $\mathbf{u}_1^T$

$$\mathbf{u}_1^T S \mathbf{u}_1 = \lambda_1 \quad (25)$$

Derive by  $\mathbf{u}_1$

We get

$$S\mathbf{u}_1 = \lambda_1\mathbf{u}_1 \quad (24)$$

Then

$\mathbf{u}_1$  is an eigenvector of  $S$ .

If we left-multiply by  $\mathbf{u}_1^T$

$$\mathbf{u}_1^T S \mathbf{u}_1 = \lambda_1 \quad (25)$$

Derive by  $\mathbf{u}_1$

We get

$$S\mathbf{u}_1 = \lambda_1\mathbf{u}_1 \quad (24)$$

Then

$\mathbf{u}_1$  is an eigenvector of  $S$ .

If we left-multiply by  $\mathbf{u}_1$

$$\mathbf{u}_1^T S \mathbf{u}_1 = \lambda_1 \quad (25)$$



## What about the second eigenvector $\mathbf{u}_2$

We have the following optimization problem

$$\begin{aligned} \max \quad & \mathbf{u}_2^T S \mathbf{u}_2 \\ \text{s.t.} \quad & \mathbf{u}_2^T \mathbf{u}_2 = 1 \\ & \mathbf{u}_2^T \mathbf{u}_1 = 0 \end{aligned}$$

Lagrangian

$$L(\mathbf{u}_2, \lambda_1, \lambda_2) = \mathbf{u}_2^T S \mathbf{u}_2 - \lambda_1 (\mathbf{u}_2^T \mathbf{u}_2 - 1) - \lambda_2 (\mathbf{u}_2^T \mathbf{u}_1 - 0)$$

## What about the second eigenvector $\mathbf{u}_2$

We have the following optimization problem

$$\begin{aligned} \max \quad & \mathbf{u}_2^T S \mathbf{u}_2 \\ \text{s.t.} \quad & \mathbf{u}_2^T \mathbf{u}_2 = 1 \\ & \mathbf{u}_2^T \mathbf{u}_1 = 0 \end{aligned}$$

Lagrangian

$$L(\mathbf{u}_2, \lambda_1, \lambda_2) = \mathbf{u}_2^T S \mathbf{u}_2 - \lambda_1 (\mathbf{u}_2^T \mathbf{u}_2 - 1) - \lambda_2 (\mathbf{u}_2^T \mathbf{u}_1 - 0)$$

# Explanation

## First the constrained minimization

- We want to maximize  $\mathbf{u}_2^T \mathbf{S} \mathbf{u}_2$

Given that the second eigenvector is orthonormal

- We have then  $\mathbf{u}_2^T \mathbf{u}_2 = 1$

Under orthonormal vectors

- The covariance goes to zero

$$\text{cov}(\mathbf{u}_1, \mathbf{u}_2) = \mathbf{u}_2^T \mathbf{S} \mathbf{u}_1 = \mathbf{u}_2^T \lambda_1 \mathbf{u}_1 = \lambda_1 \mathbf{u}_1^T \mathbf{u}_2 = 0$$

# Explanation

## First the constrained minimization

- We want to maximize  $\mathbf{u}_2^T \mathbf{S} \mathbf{u}_2$

## Given that the second eigenvector is orthonormal

- We have then  $\mathbf{u}_2^T \mathbf{u}_2 = 1$

## Under orthonormal vectors

- The covariance goes to zero

$$\text{cov}(\mathbf{u}_1, \mathbf{u}_2) = \mathbf{u}_2^T \mathbf{S} \mathbf{u}_1 = \mathbf{u}_2^T \lambda_1 \mathbf{u}_1 = \lambda_1 \mathbf{u}_1^T \mathbf{u}_2 = 0$$

# Explanation

## First the constrained minimization

- We want to maximize  $\mathbf{u}_2^T S \mathbf{u}_2$

## Given that the second eigenvector is orthonormal

- We have then  $\mathbf{u}_2^T \mathbf{u}_2 = 1$

## Under orthonormal vectors

- The covariance goes to zero  
$$\text{cov}(\mathbf{u}_1, \mathbf{u}_2) = \mathbf{u}_2^T S \mathbf{u}_1 = \mathbf{u}_2^T \lambda_1 \mathbf{u}_1 = \lambda_1 \mathbf{u}_2^T \mathbf{u}_1 = 0$$

# Meaning

The PCA's are perpendicular

$$L(\mathbf{u}_2, \lambda_1, \lambda_2) = \mathbf{u}_2^T S \mathbf{u}_2 - \lambda_1 (\mathbf{u}_2^T \mathbf{u}_2 - 1) - \lambda_2 (\mathbf{u}_2^T \mathbf{u}_1 - 0)$$

Take the derivative with respect to  $\mathbf{u}_2$

$$\frac{\partial L(\mathbf{u}_2, \lambda_1, \lambda_2)}{\partial \mathbf{u}_2} = S \mathbf{u}_2 - \lambda_1 \mathbf{u}_2 - \lambda_2 \mathbf{u}_1 = 0$$

Then, we left multiply  $\mathbf{u}_1^T$

$$\mathbf{u}_1^T S \mathbf{u}_2 - \lambda_1 \mathbf{u}_1^T \mathbf{u}_2 - \lambda_2 \mathbf{u}_1^T \mathbf{u}_1 = 0$$

# Meaning

The PCA's are perpendicular

$$L(\mathbf{u}_2, \lambda_1, \lambda_2) = \mathbf{u}_2^T S \mathbf{u}_2 - \lambda_1 (\mathbf{u}_2^T \mathbf{u}_2 - 1) - \lambda_2 (\mathbf{u}_2^T \mathbf{u}_1 - 0)$$

The the derivative with respect to  $\mathbf{u}_2$

$$\frac{\partial L(\mathbf{u}_2, \lambda_1, \lambda_2)}{\partial \mathbf{u}_2} = S \mathbf{u}_2 - \lambda_1 \mathbf{u}_2 - \lambda_2 \mathbf{u}_1 = 0$$

Then, we left multiply  $\mathbf{u}_1^T$

$$\mathbf{u}_1^T S \mathbf{u}_2 - \lambda_1 \mathbf{u}_1^T \mathbf{u}_2 - \lambda_2 \mathbf{u}_1^T \mathbf{u}_1 = 0$$

# Meaning

The PCA's are perpendicular

$$L(\mathbf{u}_2, \lambda_1, \lambda_2) = \mathbf{u}_2^T S \mathbf{u}_2 - \lambda_1 (\mathbf{u}_2^T \mathbf{u}_2 - 1) - \lambda_2 (\mathbf{u}_2^T \mathbf{u}_1 - 0)$$

The the derivative with respect to  $\mathbf{u}_2$

$$\frac{\partial L(\mathbf{u}_2, \lambda_1, \lambda_2)}{\partial \mathbf{u}_2} = S \mathbf{u}_2 - \lambda_1 \mathbf{u}_2 - \lambda_2 \mathbf{u}_1 = 0$$

Then, we left multiply  $\mathbf{u}_1$

$$\mathbf{u}_1^T S \mathbf{u}_2 - \lambda_1 \mathbf{u}_1^T \mathbf{u}_2 - \lambda_2 \mathbf{u}_1^T \mathbf{u}_1 = 0$$



Then, we have that

### Something Notable

$$0 - 0 - \lambda_2 = 0$$

We have

$$S\mathbf{u}_2 - \lambda_2\mathbf{u}_2 = \mathbf{0}$$

implying

- $\mathbf{u}_2$  is the eigenvector of  $S$  with second largest eigenvalue  $\lambda_2$ .

Then, we have that

### Something Notable

$$0 - 0 - \lambda_2 = 0$$

We have

$$S\mathbf{u}_2 - \lambda_2\mathbf{u}_2 = 0$$

implying

- $\mathbf{u}_2$  is the eigenvector of  $S$  with second largest eigenvalue  $\lambda_2$ .

Then, we have that

### Something Notable

$$0 - 0 - \lambda_2 = 0$$

We have

$$S\mathbf{u}_2 - \lambda_2\mathbf{u}_2 = \mathbf{0}$$

Implying

- $\mathbf{u}_2$  is the eigenvector of  $S$  with second largest eigenvalue  $\lambda_2$ .

Thus

Variance will be the maximum when

$$\mathbf{u}_1^T S \mathbf{u}_1 = \lambda_1 \quad (26)$$

is set to the largest eigenvalue. Also known as the First Principal Component

By Induction

It is possible for  $M$ -dimensional space to define  $M$  eigenvectors  $\mathbf{u}_1, \mathbf{u}_2, \dots, \mathbf{u}_M$  of the data covariance  $S$  corresponding to  $\lambda_1, \lambda_2, \dots, \lambda_M$  that maximize the variance of the projected data.

Computational Cost

- Full eigenvector decomposition  $O(d^3)$
- Power Method  $O(Md^2)$  "Golub and Van Loan, 1996"
- Use the Expectation Maximization Algorithm

Thus

Variance will be the maximum when

$$\mathbf{u}_1^T \mathbf{S} \mathbf{u}_1 = \lambda_1 \quad (26)$$

is set to the largest eigenvalue. Also known as the First Principal Component

## By Induction

It is possible for  $M$ -dimensional space to define  $M$  eigenvectors  $\mathbf{u}_1, \mathbf{u}_2, \dots, \mathbf{u}_M$  of the data covariance  $\mathbf{S}$  corresponding to  $\lambda_1, \lambda_2, \dots, \lambda_M$  that maximize the variance of the projected data.

## Computational Costs

- Full eigenvector decomposition  $O(d^3)$
- Power Method  $O(Md^2)$  "Golub and Van Loan, 1996"
- Use the Expectation Maximization Algorithm

Thus

Variance will be the maximum when

$$\mathbf{u}_1^T \mathbf{S} \mathbf{u}_1 = \lambda_1 \quad (26)$$

is set to the largest eigenvalue. Also known as the First Principal Component

## By Induction

It is possible for  $M$ -dimensional space to define  $M$  eigenvectors  $\mathbf{u}_1, \mathbf{u}_2, \dots, \mathbf{u}_M$  of the data covariance  $\mathbf{S}$  corresponding to  $\lambda_1, \lambda_2, \dots, \lambda_M$  that maximize the variance of the projected data.

## Computational Cost

- 1 Full eigenvector decomposition  $O(d^3)$
- 2 Power Method  $O(Md^2)$  “Golub and Van Loan, 1996”
- 3 Use the Expectation Maximization Algorithm

# Outline

## 1 Fisher Linear Discriminant

- Introduction
- The Rotation Idea
- Solution
  - Scatter measure
- The Cost Function

## 2 Principal Components and Singular Value Decomposition

- Introduction
- **Principal Component Analysis AKA Karhunen-Loeve Transform**
  - Projecting the Data
  - Lagrange Multipliers
  - **The Process**
  - Example
- Singular Value Decomposition
  - Introduction
  - Image Compression

We have the following steps

Determine covariance matrix

$$S = \frac{1}{N-1} \sum_{i=1}^N (\mathbf{x}_i - \bar{\mathbf{x}}) (\mathbf{x}_i - \bar{\mathbf{x}})^T \quad (27)$$

Generate the decomposition

$$S = U \Sigma U^T$$

With

- Eigenvalues in  $\Sigma$  and eigenvectors in the columns of  $U$ .



We have the following steps

Determine covariance matrix

$$S = \frac{1}{N-1} \sum_{i=1}^N (\mathbf{x}_i - \bar{\mathbf{x}}) (\mathbf{x}_i - \bar{\mathbf{x}})^T \quad (27)$$

Generate the decomposition

$$S = U \Sigma U^T$$

With

- Eigenvalues in  $\Sigma$  and eigenvectors in the columns of  $U$ .

We have the following steps

Determine covariance matrix

$$S = \frac{1}{N-1} \sum_{i=1}^N (\mathbf{x}_i - \bar{\mathbf{x}}) (\mathbf{x}_i - \bar{\mathbf{x}})^T \quad (27)$$

Generate the decomposition

$$S = U \Sigma U^T$$

With

- Eigenvalues in  $\Sigma$  and eigenvectors in the columns of  $U$ .

Then

Project samples  $\mathbf{x}_i$  into subspaces  $\text{dim}=k$

$$\mathbf{z}_i = \mathbf{U}_K^T \mathbf{x}_i$$

- With  $\mathbf{U}_k$  is a matrix with  $k$  columns

# Outline

## 1 Fisher Linear Discriminant

- Introduction
- The Rotation Idea
- Solution
  - Scatter measure
- The Cost Function

## 2 Principal Components and Singular Value Decomposition

- Introduction
- **Principal Component Analysis AKA Karhunen-Loeve Transform**
  - Projecting the Data
  - Lagrange Multipliers
  - The Process
  - **Example**
- Singular Value Decomposition
  - Introduction
  - Image Compression

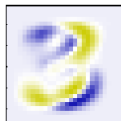
# Example

From Bishop

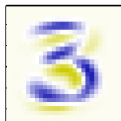
Mean



$\lambda_1 = 3.4 \cdot 10^5$



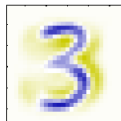
$\lambda_2 = 2.8 \cdot 10^5$



$\lambda_3 = 2.4 \cdot 10^5$

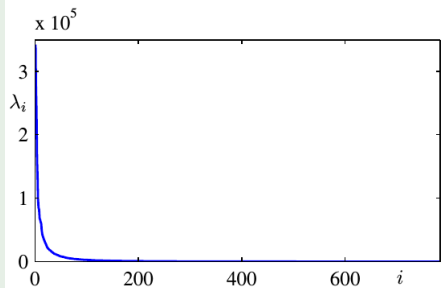


$\lambda_4 = 1.6 \cdot 10^5$



# Example

From Bishop



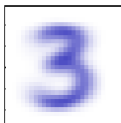
# Example

From Bishop

Original



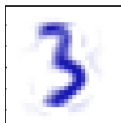
$M = 1$



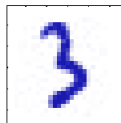
$M = 10$



$M = 50$



$M = 250$



# Outline

## 1 Fisher Linear Discriminant

- Introduction
- The Rotation Idea
- Solution
  - Scatter measure
- The Cost Function

## 2 Principal Components and Singular Value Decomposition

- Introduction
- Principal Component Analysis AKA Karhunen-Loeve Transform
  - Projecting the Data
  - Lagrange Multipliers
  - The Process
  - Example
- **Singular Value Decomposition**
  - Introduction
  - Image Compression



# Outline

## 1 Fisher Linear Discriminant

- Introduction
- The Rotation Idea
- Solution
  - Scatter measure
- The Cost Function

## 2 Principal Components and Singular Value Decomposition

- Introduction
- Principal Component Analysis AKA Karhunen-Loeve Transform
  - Projecting the Data
  - Lagrange Multipliers
  - The Process
  - Example
- Singular Value Decomposition
  - Introduction
  - Image Compression

# What happened with no-square matrices

We can still diagonalize it

Thus, we can obtain certain properties.

We want to avoid the problems with

$$S^{-1}AS$$

The eigenvectors in  $S$  have three big problems.

- They are usually not orthogonal.
- There are not always enough eigenvectors.
- $Ax = \lambda x$  requires  $A$  to be square.

# What happened with no-square matrices

We can still diagonalize it

Thus, we can obtain certain properties.

We want to avoid the problems with

$$S^{-1}AS$$

The eigenvectors in  $S$  have three big problems.

- They are usually not orthogonal.
- There are not always enough eigenvectors.
- $Ax = \lambda x$  requires  $A$  to be square.

# What happened with no-square matrices

We can still diagonalize it

Thus, we can obtain certain properties.

We want to avoid the problems with

$$S^{-1}AS$$

The eigenvectors in  $S$  have three big problems

- 1 They are usually not orthogonal.
- 2 There are not always enough eigenvectors.
- 3  $Ax = \lambda x$  requires  $A$  to be square.

Therefore, we can look at the following problem

We have a series of vectors

$$\{\mathbf{x}_1, \mathbf{x}_2, \dots, \mathbf{x}_d\}$$

Then imagine a set of projection vectors and differences

$$\{\beta_1, \beta_2, \dots, \beta_d\} \text{ and } \{\alpha_1, \alpha_2, \dots, \alpha_d\}$$

We want to know a little bit of the relations between them

- After all, we are looking at the possibility of using them for our problem

Therefore, we can look at the following problem

We have a series of vectors

$$\{\mathbf{x}_1, \mathbf{x}_2, \dots, \mathbf{x}_d\}$$

Then imagine a set of projection vectors and differences

$$\{\beta_1, \beta_2, \dots, \beta_d\} \text{ and } \{\alpha_1, \alpha_2, \dots, \alpha_d\}$$

We want to know a little bit of the relations between them

- After all, we are looking at the possibility of using them for our problem

Therefore, we can look at the following problem

We have a series of vectors

$$\{\mathbf{x}_1, \mathbf{x}_2, \dots, \mathbf{x}_d\}$$

Then imagine a set of projection vectors and differences

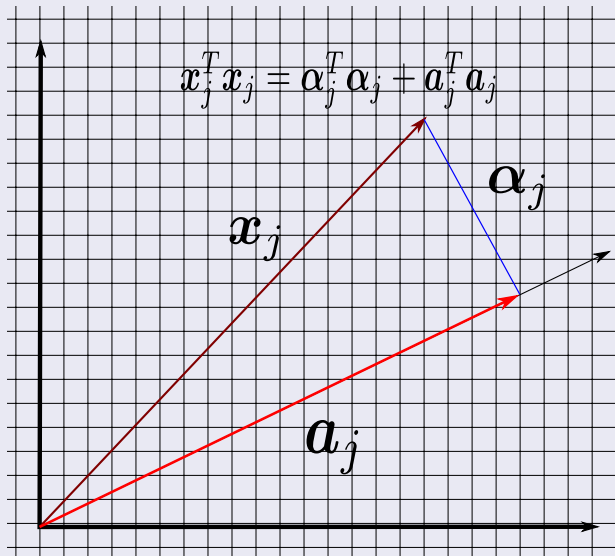
$$\{\beta_1, \beta_2, \dots, \beta_d\} \text{ and } \{\alpha_1, \alpha_2, \dots, \alpha_d\}$$

We want to know a little bit of the relations between them

- After all, we are looking at the possibility of using them for our problem

# Using the Hypotenuse

A little bit of Geometry, we get





Therefore

We have two possible quantities for each  $j$

$$\alpha_j^T \alpha_j = \mathbf{x}_j^T \mathbf{x}_j - \mathbf{a}_j^T \mathbf{a}_j$$

$$\mathbf{a}_j^T \mathbf{a}_j = \mathbf{x}_j^T \mathbf{x}_j - \alpha_j^T \alpha_j$$

Then we can minimize and maximize given that  $\mathbf{x}_j^T \mathbf{x}_j$  is a constant

$$\min \sum_{j=1}^n \alpha_j^T \alpha_j$$

$$\max \sum_{j=1}^n \mathbf{a}_j^T \mathbf{a}_j$$

Therefore

We have two possible quantities for each  $j$

$$\alpha_j^T \alpha_j = x_j^T x_j - a_j^T a_j$$

$$a_j^T a_j = x_j^T x_j - \alpha_j^T \alpha_j$$

Then, we can minimize and maximize given that  $x_j^T x_j$  is a constant

$$\min \sum_{j=1}^n \alpha_j^T \alpha_j$$

$$\max \sum_{j=1}^n a_j^T a_j$$

# Actually this is known as the dual problem (Weak Duality)

An example of this

$$\begin{aligned} \min \quad & \mathbf{w}^T \mathbf{x} \\ \text{s.t.} \quad & \mathbf{Ax} \leq \mathbf{b} \\ & \mathbf{x} \geq 0 \end{aligned}$$

Then, using what is known as slack variables

$$\mathbf{Ax} + \mathbf{A}'\mathbf{x}' = \mathbf{b}$$

Each row lives in the column space, but the  $y_i$  lives in the column space

$$(\mathbf{Ax} + \mathbf{A}'\mathbf{x}')_i \rightarrow y_i \text{ and } \mathbf{x}' \geq 0$$

Actually this is known as the dual problem (Weak Duality)

An example of this

$$\begin{aligned} \min \quad & \mathbf{w}^T \mathbf{x} \\ \text{s.t.} \quad & \mathbf{Ax} \leq \mathbf{b} \\ & \mathbf{x} \geq 0 \end{aligned}$$

Then, using what is known as slack variables

$$\mathbf{Ax} + \mathbf{A}'\mathbf{x}' = \mathbf{b}$$

Each row lives in the column space, but the  $y_i$  lives in the column space

$$(\mathbf{Ax} + \mathbf{A}'\mathbf{x}')_i \rightarrow y_i \text{ and } \mathbf{x}' \geq 0$$

Actually this is known as the dual problem (Weak Duality)

An example of this

$$\begin{aligned} \min \quad & \mathbf{w}^T \mathbf{x} \\ \text{s.t.} \quad & \mathbf{Ax} \leq \mathbf{b} \\ & \mathbf{x} \geq 0 \end{aligned}$$

Then, using what is known as slack variables

$$\mathbf{Ax} + \mathbf{A}'\mathbf{x} = \mathbf{b}$$

Each row lives in the column space, but the  $y_i$  lives in the column space

$$(\mathbf{Ax} + \mathbf{A}'\mathbf{x})_i \rightarrow y_i \text{ and } \mathbf{x}' \geq 0$$

Then, we have that

### Example

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

Element in the column space of dimensionality have three dimensions

- But in the row space their dimension is 2

### Properties

Then, we have that

### Example

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

Element in the column space of dimensionality have three dimensions

- But in the row space their dimension is 2

### Properties

Then, we have that

### Example

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

Element in the column space of dimensionality have three dimensions

- But in the row space their dimension is 2

### Properties



We have then

Stack such vectors that in the  $d$ -dimensional space

- In a matrix  $A$  of  $n \times d$

$$A = \begin{bmatrix} \mathbf{a}_1^T \\ \mathbf{a}_2^T \\ \vdots \\ \mathbf{a}_n^T \end{bmatrix}$$

The matrix works as a Projection Matrix

- We are looking for a unit vector  $\mathbf{v}$  such that length of the projection is maximized.

We have then

Stack such vectors that in the  $d$ -dimensional space

- In a matrix  $A$  of  $n \times d$

$$A = \begin{bmatrix} \mathbf{a}_1^T \\ \mathbf{a}_2^T \\ \vdots \\ \mathbf{a}_n^T \end{bmatrix}$$

The matrix works as a Projection Matrix

- We are looking for a unit vector  $\mathbf{v}$  such that length of the projection is maximized.

Why? Do you remember the Projection to a single vector  $p$ ?

Definition of the projection under unitary vector

$$p = \frac{v^T a_i}{v^T v} v = \left[ v^T a_i \right] v$$

Therefore the length of the projected vector is

$$\left\| \left[ v^T a_i \right] v \right\| = \left| v^T a_i \right|$$

Why? Do you remember the Projection to a single vector  $p$ ?

Definition of the projection under unitary vector

$$p = \frac{v^T a_i}{v^T v} v = \left[ v^T a_i \right] v$$

Therefore the length of the projected vector is

$$\left\| \left[ v^T a_i \right] v \right\| = \left| v^T a_i \right|$$

Then

Thus with a little bit of notation

$$A\mathbf{v} = \begin{bmatrix} \mathbf{a}_1^T \\ \mathbf{a}_2^T \\ \vdots \\ \mathbf{a}_d^T \end{bmatrix} \mathbf{v} = \begin{bmatrix} \mathbf{a}_1^T \mathbf{v} \\ \mathbf{a}_2^T \mathbf{v} \\ \vdots \\ \mathbf{a}_d^T \mathbf{v} \end{bmatrix}$$

Therefore

$$\|A\mathbf{v}\| = \sqrt{\sum_{i=1}^d (\mathbf{a}_i^T \mathbf{v})^2}$$

Then

Thus with a little bit of notation

$$A\mathbf{v} = \begin{bmatrix} \mathbf{a}_1^T \\ \mathbf{a}_2^T \\ \vdots \\ \mathbf{a}_d^T \end{bmatrix} \mathbf{v} = \begin{bmatrix} \mathbf{a}_1^T \mathbf{v} \\ \mathbf{a}_2^T \mathbf{v} \\ \vdots \\ \mathbf{a}_d^T \mathbf{v} \end{bmatrix}$$

Therefore

$$\|A\mathbf{v}\| = \sqrt{\sum_{i=1}^d (\mathbf{a}_i^T \mathbf{v})^2}$$

Then

It is possible to ask to maximize the longitude of such vector  
(Singular Vector)

$$\mathbf{v}_1 = \arg \max_{\|\mathbf{v}\|=1} \|A\mathbf{v}\|$$

Then, we can define the following singular value

$$\sigma_1(A) = \|A\mathbf{v}_1\|$$

Then

It is possible to ask to maximize the longitude of such vector  
(Singular Vector)

$$\mathbf{v}_1 = \arg \max_{\|\mathbf{v}\|=1} \|\mathbf{A}\mathbf{v}\|$$

Then, we can define the following singular value

$$\sigma_1(A) = \|\mathbf{A}\mathbf{v}_1\|$$



# This is known as

## Definition

- The **best-fit line problem** describes the problem of finding the best line for a set of data points, where the quality of the line is measured by the sum of squared (perpendicular) distances of the points to the line.
  - ▶ Remember, we are looking at the dual problem....

## Generalization

- This can be transferred to higher dimensions: One can find the best-fit  $d$ -dimensional subspace, so the subspace which minimizes the sum of the squared distances of the points to the subspace

# This is known as

## Definition

- The **best-fit line problem** describes the problem of finding the best line for a set of data points, where the quality of the line is measured by the sum of squared (perpendicular) distances of the points to the line.
  - ▶ Remember, we are looking at the dual problem....

## Generalization

- This can be transferred to higher dimensions: One can find the best-fit  $d$ -dimensional subspace, so the subspace which minimizes the sum of the squared distances of the points to the subspace

## Then, in a Greedy Fashion

The second singular vector  $v_2$

$$v_2 = \arg \max_{v \perp v_1, \|v\|=1} \|Av\|$$

Then you go through this process:

- Stop when we have found all the following vectors:

$$v_1, v_2, \dots, v_r$$

As singular vectors and

$$\arg \max_{\substack{v \perp v_1, v_2, \dots, v_r \\ \|v\|=1}} \|Av\|$$

## Then, in a Greedy Fashion

The second singular vector  $v_2$

$$v_2 = \arg \max_{v \perp v_1, \|v\|=1} \|Av\|$$

Then you go through this process

- Stop when we have found all the following vectors:

$$v_1, v_2, \dots, v_r$$

As singular vectors and

$$\arg \max_{\substack{v \perp v_1, v_2, \dots, v_r \\ \|v\|=1}} \|Av\|$$

## Then, in a Greedy Fashion

The second singular vector  $v_2$

$$v_2 = \arg \max_{v \perp v_1, \|v\|=1} \|Av\|$$

Then you go through this process

- Stop when we have found all the following vectors:

$$v_1, v_2, \dots, v_r$$

As singular vectors and

$$\arg \max_{\substack{v \perp v_1, v_2, \dots, v_r \\ \|v\| = 1}} \|Av\|$$

# Proving that the strategy is good

## Theorem

- Let  $A$  be an  $n \times d$  matrix where  $\mathbf{v}_1, \mathbf{v}_2, \dots, \mathbf{v}_r$  are the singular vectors defined above. For  $1 \leq k \leq r$ , let  $V_k$  be the subspace spanned by  $\mathbf{v}_1, \mathbf{v}_2, \dots, \mathbf{v}_k$ . Then for each  $k$ ,  $V_k$  is the best-fit  $k$ -dimensional subspace for  $A$ .

# Proof

For  $k = 1$

- What about  $k = 2$ ? Let  $W$  be a best-fit 2- dimensional subspace for  $A$ .

For any basis  $w_1, w_2$  of  $W$

- $|Aw_1|^2 + |Aw_2|^2$  is the sum of the squared lengths of the projections of the rows of  $A$  to  $W$ .

Now choose a basis  $w_1, w_2$  so that  $w_1$  is perpendicular to  $w_2$ .

- This can be a unit vector perpendicular to  $v_1$  projection in  $W$ .

# Proof

For  $k = 1$

- What about  $k = 2$ ? Let  $W$  be a best-fit 2- dimensional subspace for  $A$ .

For any basis  $w_1, w_2$  of  $W$

- $|Aw_1|^2 + |Aw_2|^2$  is the sum of the squared lengths of the projections of the rows of  $A$  to  $W$ .

Now, choose a basis  $w_1, w_2$  so that  $w_2$  is perpendicular to  $w_1$ .

- This can be a unit vector perpendicular to  $w_1$  projection in  $W$ .



# Proof

For  $k = 1$

- What about  $k = 2$ ? Let  $W$  be a best-fit 2- dimensional subspace for  $A$ .

For any basis  $w_1, w_2$  of  $W$

- $|Aw_1|^2 + |Aw_2|^2$  is the sum of the squared lengths of the projections of the rows of  $A$  to  $W$ .

Now, choose a basis  $w_1, w_2$  so that  $w_2$  is perpendicular to  $v_1$

- This can be a unit vector perpendicular to  $v_1$  projection in  $W$ .

Do you remember  $\mathbf{v}_1 = \arg \max_{\|\mathbf{v}\|=1} \|A\mathbf{v}\|$ ?

Therefore

$$|A\mathbf{w}_1|^2 \leq |A\mathbf{v}_1|^2 \text{ and } |A\mathbf{w}_2|^2 \leq |A\mathbf{v}_2|^2$$

Then

$$|A\mathbf{w}_1|^2 + |A\mathbf{w}_2|^2 \leq |A\mathbf{v}_1|^2 + |A\mathbf{v}_2|^2$$

In a similar way, for  $k=2, \dots, n$

- $V_k$  is at least as good as  $W$  and hence is optimal.

Do you remember  $\mathbf{v}_1 = \arg \max_{\|\mathbf{v}\|=1} \|A\mathbf{v}\|$ ?

Therefore

$$|A\mathbf{w}_1|^2 \leq |A\mathbf{v}_1|^2 \text{ and } |A\mathbf{w}_2|^2 \leq |A\mathbf{v}_2|^2$$

Then

$$|A\mathbf{w}_1|^2 + |A\mathbf{w}_2|^2 \leq |A\mathbf{v}_1|^2 + |A\mathbf{v}_2|^2$$

In a similar way, we get

- $V_k$  is at least as good as  $W$  and hence is optimal.

Do you remember  $\mathbf{v}_1 = \arg \max_{\|\mathbf{v}\|=1} \|A\mathbf{v}\|$ ?

Therefore

$$|A\mathbf{w}_1|^2 \leq |A\mathbf{v}_1|^2 \text{ and } |A\mathbf{w}_2|^2 \leq |A\mathbf{v}_2|^2$$

Then

$$|A\mathbf{w}_1|^2 + |A\mathbf{w}_2|^2 \leq |A\mathbf{v}_1|^2 + |A\mathbf{v}_2|^2$$

In a similar way for  $k > 2$

- $V_k$  is at least as good as  $W$  and hence is optimal.

## Remarks

Every Matrix has a singular value decomposition

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

Where

- The columns of  $U$  are an orthonormal basis for the column space.

Every Matrix has a singular value decomposition

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

Where

- The columns of  $U$  are an orthonormal basis for the column space.
- The columns of  $V$  are an orthonormal basis for the row space.
- The  $\Sigma$  is diagonal and the entries on its diagonal  $\sigma_i = \Sigma_{ii}$  are positive real numbers, called the singular values of  $A$ .

Every Matrix has a singular value decomposition

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

Where

- The columns of  $U$  are an orthonormal basis for the column space.
- The columns of  $V$  are an orthonormal basis for the row space.
- The  $\Sigma$  is diagonal and the entries on its diagonal  $\sigma_i = \Sigma_{ii}$  are positive real numbers, called the singular values of  $A$ .

# Properties of the Singular Value Decomposition

## First

The eigenvalues of the symmetric matrix  $A^T A$  are equal to the square of the singular values of  $A$

$$A^T A = V \Sigma U^T U^T \Sigma V^T = V \Sigma^2 V^T$$

## Second

The rank of a matrix is equal to the number of non-zero singular values.



# Properties of the Singular Value Decomposition

## First

The eigenvalues of the symmetric matrix  $A^T A$  are equal to the square of the singular values of  $A$

$$A^T A = V \Sigma U^T U^T \Sigma V^T = V \Sigma^2 V^T$$

## Second

The rank of a matrix is equal to the number of non-zero singular values.

# Outline

## 1 Fisher Linear Discriminant

- Introduction
- The Rotation Idea
- Solution
  - Scatter measure
- The Cost Function

## 2 Principal Components and Singular Value Decomposition

- Introduction
- Principal Component Analysis AKA Karhunen-Loeve Transform
  - Projecting the Data
  - Lagrange Multipliers
  - The Process
  - Example
- Singular Value Decomposition
  - Introduction
  - Image Compression

# Singular Value Decomposition as Sums

The singular value decomposition can be viewed as a sum of rank 1 matrices

$$A = A_1 + A_2 + \dots + A_R \quad (28)$$

Why?

$$\begin{aligned} u_1 A = U \begin{pmatrix} \sigma_1 & 0 & \dots & 0 \\ 0 & \sigma_2 & \dots & 0 \\ \vdots & \vdots & \ddots & \vdots \\ 0 & 0 & \dots & \sigma_R \end{pmatrix} V^T = \begin{pmatrix} u_1 & u_2 & \dots & u_R \end{pmatrix} \begin{pmatrix} \sigma_1 v_1^T \\ \sigma_2 v_2^T \\ \vdots \\ \sigma_R v_R^T \end{pmatrix} \\ = \sigma_1 u_1 v_1^T + \sigma_2 u_2 v_2^T + \dots + \sigma_R u_R v_R^T \end{aligned}$$

# Singular Value Decomposition as Sums

The singular value decomposition can be viewed as a sum of rank 1 matrices

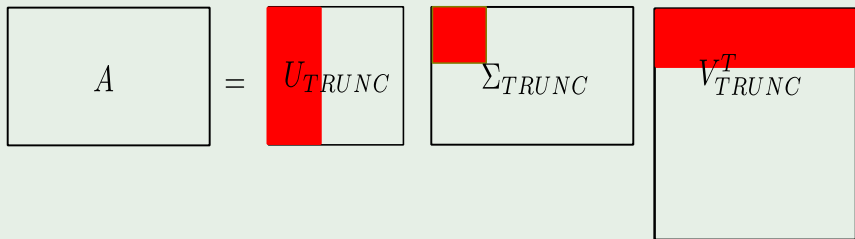
$$A = A_1 + A_2 + \dots + A_R \quad (28)$$

Why?

$$\begin{aligned} \mathbf{u}_1 A = U \begin{pmatrix} \sigma_1 & 0 & \cdots & 0 \\ 0 & \sigma_2 & \cdots & 0 \\ \vdots & \vdots & \ddots & \vdots \\ 0 & 0 & \cdots & \sigma_R \end{pmatrix} V^T = \begin{pmatrix} \mathbf{u}_1 & \mathbf{u}_2 & \cdots & \mathbf{u}_R \end{pmatrix} \begin{pmatrix} \sigma_1 \mathbf{v}_1^T \\ \sigma_2 \mathbf{v}_2^T \\ \vdots \\ \sigma_R \mathbf{v}_R^T \end{pmatrix} \\ = \sigma_1 \mathbf{u}_1 \mathbf{v}_1^T + \sigma_2 \mathbf{u}_2 \mathbf{v}_2^T + \cdots + \sigma_R \mathbf{u}_R \mathbf{v}_R^T \end{aligned}$$

## Truncating

Truncating the singular value decomposition allows us to represent the matrix with less parameters

$$A = U_{TRUNC} \Sigma_{TRUNC} V_{TRUNC}^T$$


# Truncating

Truncating the singular value decomposition allows us to represent the matrix with less parameters

$$A = U_{TRUNC} \Sigma_{TRUNC} V_{TRUNC}^T$$

For a  $512 \times 512$

- Full Representation  $512 \times 512 = 262,144$
- Rank 10 approximation  $512 \times 10 + 10 + 10 \times 512 = 10,250$
- Rank 40 approximation  $512 \times 40 + 40 + 40 \times 512 = 41,000$
- Rank 80 approximation  $512 \times 80 + 80 + 80 \times 512 = 82,000$

# Truncating

Truncating the singular value decomposition allows us to represent the matrix with less parameters

$$A = U_{TRUNC} \Sigma_{TRUNC} V_{TRUNC}^T$$

For a  $512 \times 512$

- Full Representation  $512 \times 512 = 262,144$
- Rank 10 approximation  $512 \times 10 + 10 + 10 \times 512 = 10,250$
- Rank 40 approximation  $512 \times 40 + 40 + 40 \times 512 = 41,000$
- Rank 80 approximation  $512 \times 80 + 80 + 80 \times 512 = 82,000$

# Truncating

Truncating the singular value decomposition allows us to represent the matrix with less parameters

$$A = U_{TRUNC} \Sigma_{TRUNC} V_{TRUNC}^T$$

For a  $512 \times 512$

- Full Representation  $512 \times 512 = 262,144$
- Rank 10 approximation  $512 \times 10 + 10 + 10 \times 512 = 10,250$
- Rank 40 approximation  $512 \times 40 + 40 + 40 \times 512 = 41,000$
- Rank 80 approximation  $512 \times 80 + 80 + 80 \times 512 = 82,000$