

Mathematics for Artificial Intelligence

Transformation and Applications

Andres Mendez-Vazquez

April 9, 2020

Outline

1 Linear Transformation

- Introduction
- Functions that can be defined using matrices
- Linear Functions
- Kernel and Range
- The Matrix of a Linear Transformation
- Going Back to Homogeneous Equations
- The Rank-Nullity Theorem

2 Derivative of Transformations

- Introduction
- Derivative of a Linear Transformation
- Derivative of a Quadratic Transformation

3 Linear Regression

- The Simplest Functions
- Splitting the Space
- Defining the Decision Surface
- Properties of the Hyperplane $w^T x + w_0$
- Augmenting the Vector
- Least Squared Error Procedure
 - The Geometry of a Two-Category Linearly-Separable Case
- The Error Idea
- The Final Error Equation

4 Principal Component Analysis

- Karhunen-Loeve Transform
- Projecting the Data
- Lagrange Multipliers
- The Process
- Example

5 Singular Value Decomposition

- Introduction
- Image Compression



Outline

1 Linear Transformation

● Introduction

- Functions that can be defined using matrices
- Linear Functions
- Kernel and Range
- The Matrix of a Linear Transformation
- Going Back to Homogeneous Equations
- The Rank-Nullity Theorem

2 Derivative of Transformations

- Introduction
- Derivative of a Linear Transformation
- Derivative of a Quadratic Transformation

3 Linear Regression

- The Simplest Functions
- Splitting the Space
- Defining the Decision Surface
- Properties of the Hyperplane $w^T x + w_0$
- Augmenting the Vector
- Least Squared Error Procedure
 - The Geometry of a Two-Category Linearly-Separable Case
- The Error Idea
- The Final Error Equation

4 Principal Component Analysis

- Karhunen-Loeve Transform
- Projecting the Data
- Lagrange Multipliers
- The Process
- Example

5 Singular Value Decomposition

- Introduction
- Image Compression



Going further than solving $Ax = y$

We can go further

We can think on the matrix A as a function!!!

Going further than solving $A\mathbf{x} = \mathbf{y}$

We can go further

We can think on the matrix A as a function!!!

In general

A function f whose domain \mathbb{R}^n and defines a rule that associate $\mathbf{x} \in \mathbb{R}^n$ to a vector $\mathbf{y} \in \mathbb{R}^m$

$$\mathbf{y} = f(\mathbf{x}) \text{ equivalently } f : \mathbb{R}^n \longrightarrow \mathbb{R}^m$$

Going further than solving $Ax = y$

We can go further

We can think on the matrix A as a function!!!

In general

A function f whose domain \mathbb{R}^n and defines a rule that associate $x \in \mathbb{R}^n$ to a vector $y \in \mathbb{R}^m$

$$y = f(x) \text{ equivalently } f : \mathbb{R}^n \longrightarrow \mathbb{R}^m$$

We like the second expression because

① It is easy to identify the domain \mathbb{R}^n

② It is easy to find the range \mathbb{R}^m

Going further than solving $Ax = y$

We can go further

We can think on the matrix A as a function!!!

In general

A function f whose domain \mathbb{R}^n and defines a rule that associate $x \in \mathbb{R}^n$ to a vector $y \in \mathbb{R}^m$

$$y = f(x) \text{ equivalently } f : \mathbb{R}^n \longrightarrow \mathbb{R}^m$$

We like the second expression because

- 1 It is easy to identify the domain \mathbb{R}^n
- 2 It is easy to find the range \mathbb{R}^m

Examples

$$f : \mathbb{R} \rightarrow \mathbb{R}^3$$

$$f(t) = \begin{pmatrix} x(t) \\ y(t) \\ z(t) \end{pmatrix} = \begin{pmatrix} t \\ 3t^2 + 1 \\ \sin(t) \end{pmatrix}$$

This are called parametric functions

- Depending on the context, it could represent the position or the velocity of a mass point.



Examples

$$f : \mathbb{R} \rightarrow \mathbb{R}^3$$

$$f(t) = \begin{pmatrix} x(t) \\ y(t) \\ z(t) \end{pmatrix} = \begin{pmatrix} t \\ 3t^2 + 1 \\ \sin(t) \end{pmatrix}$$

This are called parametric functions

- Depending on the context, it could represent the position or the velocity of a mass point.



Outline

1 Linear Transformation

- Introduction
- **Functions that can be defined using matrices**
- Linear Functions
- Kernel and Range
- The Matrix of a Linear Transformation
- Going Back to Homogeneous Equations
- The Rank-Nullity Theorem

2 Derivative of Transformations

- Introduction
- Derivative of a Linear Transformation
- Derivative of a Quadratic Transformation

3 Linear Regression

- The Simplest Functions
- Splitting the Space
- Defining the Decision Surface
- Properties of the Hyperplane $w^T x + w_0$
- Augmenting the Vector
- Least Squared Error Procedure
 - The Geometry of a Two-Category Linearly-Separable Case
- The Error Idea
- The Final Error Equation

4 Principal Component Analysis

- Karhunen-Loeve Transform
- Projecting the Data
- Lagrange Multipliers
- The Process
- Example

5 Singular Value Decomposition

- Introduction
- Image Compression



A Classic Example

We have

if A is a $m \times n$, we can use A to define a function.

We will call them

$$f_A: \mathbb{R}^n \rightarrow \mathbb{R}^m$$

In other words

$$f_A(\mathbf{x}) = A\mathbf{x}$$



A Classic Example

We have

if A is a $m \times n$, we can use A to define a function.

We will call them

$$f_A : \mathbb{R}^n \rightarrow \mathbb{R}^m$$

In other words

$$f_A(\mathbf{x}) = A\mathbf{x}$$



A Classic Example

We have

if A is a $m \times n$, we can use A to define a function.

We will call them

$$f_A : \mathbb{R}^n \rightarrow \mathbb{R}^m$$

In other words

$$f_A(\mathbf{x}) = Ax$$



Example

Let

$$A_{2 \times 3} = \begin{pmatrix} 1 & 2 & 3 \\ 4 & 5 & 6 \end{pmatrix}$$

This allows to define

$$f_A \begin{pmatrix} x \\ y \\ z \end{pmatrix} = \begin{pmatrix} 1 & 2 & 3 \\ 4 & 5 & 6 \end{pmatrix} \begin{pmatrix} x \\ y \\ z \end{pmatrix} = \begin{pmatrix} x + 2y + z \\ 4x + 5y + 6z \end{pmatrix}$$

We have

- For each vector $x \in \mathbb{R}^3$ to the vector $Ax \in \mathbb{R}^2$



Example

Let

$$A_{2 \times 3} = \begin{pmatrix} 1 & 2 & 3 \\ 4 & 5 & 6 \end{pmatrix}$$

This allows to define

$$f_A \begin{pmatrix} x \\ y \\ z \end{pmatrix} = \begin{pmatrix} 1 & 2 & 3 \\ 4 & 5 & 6 \end{pmatrix} \begin{pmatrix} x \\ y \\ z \end{pmatrix} = \begin{pmatrix} x + 2y + z \\ 4x + 5y + 6z \end{pmatrix}$$

We have

- For each vector $x \in \mathbb{R}^3$ to the vector $Ax \in \mathbb{R}^2$



Example

Let

$$A_{2 \times 3} = \begin{pmatrix} 1 & 2 & 3 \\ 4 & 5 & 6 \end{pmatrix}$$

This allows to define

$$f_A \begin{pmatrix} x \\ y \\ z \end{pmatrix} = \begin{pmatrix} 1 & 2 & 3 \\ 4 & 5 & 6 \end{pmatrix} \begin{pmatrix} x \\ y \\ z \end{pmatrix} = \begin{pmatrix} x + 2y + z \\ 4x + 5y + 6z \end{pmatrix}$$

We have

- For each vector $x \in \mathbb{R}^3$ to the vector $Ax \in \mathbb{R}^2$



Outline

1 Linear Transformation

- Introduction
- Functions that can be defined using matrices
- **Linear Functions**
- Kernel and Range
- The Matrix of a Linear Transformation
- Going Back to Homogeneous Equations
- The Rank-Nullity Theorem

2 Derivative of Transformations

- Introduction
- Derivative of a Linear Transformation
- Derivative of a Quadratic Transformation

3 Linear Regression

- The Simplest Functions
- Splitting the Space
- Defining the Decision Surface
- Properties of the Hyperplane $w^T x + w_0$
- Augmenting the Vector
- Least Squared Error Procedure
 - The Geometry of a Two-Category Linearly-Separable Case
- The Error Idea
- The Final Error Equation

4 Principal Component Analysis

- Karhunen-Loeve Transform
- Projecting the Data
- Lagrange Multipliers
- The Process
- Example

5 Singular Value Decomposition

- Introduction
- Image Compression



We have

Definition

A function $f : \mathbb{R}^n \rightarrow \mathbb{R}^m$ is said to be linear if

- 1 $f(\mathbf{x}_1 + \mathbf{x}_2) = f(\mathbf{x}_1) + f(\mathbf{x}_2)$
- 2 $f(c\mathbf{x}) = cf(\mathbf{x})$

for all $\mathbf{x}_1, \mathbf{x}_2 \in \mathbb{R}^n$ and for all the scalars c .

Notes

A linear function f is also known as a linear transformation.



Cinvestav

We have

Definition

A function $f : \mathbb{R}^n \rightarrow \mathbb{R}^m$ is said to be linear if

- 1 $f(\mathbf{x}_1 + \mathbf{x}_2) = f(\mathbf{x}_1) + f(\mathbf{x}_2)$
- 2 $f(c\mathbf{x}) = cf(\mathbf{x})$

for all $\mathbf{x}_1, \mathbf{x}_2 \in \mathbb{R}^n$ and for all the scalars c .

Thus

A linear function f is also known as a linear transformation.



We have the following proposition

Proposition

$f : \mathbb{R}^n \rightarrow \mathbb{R}^m$ is linear \iff for all vectors $\mathbf{x}_1, \mathbf{x}_2 \in \mathbb{R}^n$ and for all the scalars c_1, c_2 :

$$f(c_1\mathbf{x}_1 + c_2\mathbf{x}_2) = c_1f(\mathbf{x}_1) + c_2f(\mathbf{x}_2)$$

Proof

Any idea?



We have the following proposition

Proposition

$f : \mathbb{R}^n \rightarrow \mathbb{R}^m$ is linear \iff for all vectors $\mathbf{x}_1, \mathbf{x}_2 \in \mathbb{R}^n$ and for all the scalars c_1, c_2 :

$$f(c_1\mathbf{x}_1 + c_2\mathbf{x}_2) = c_1f(\mathbf{x}_1) + c_2f(\mathbf{x}_2)$$

Proof

Any idea?



Proof

If $A_{m \times n}$ is a matrix, f_A is a linear transformation

How?

First

$$f_A(x_1 + x_2) = A(x_1 + x_2) = Ax_1 + Ax_2 = f_A(x_1) + f_A(x_2)$$

Second

What about $f_A(cx_1)$?



Proof

If $A_{m \times n}$ is a matrix, f_A is a linear transformation

How?

First

$$f_A(\mathbf{x}_1 + \mathbf{x}_2) = A(\mathbf{x}_1 + \mathbf{x}_2) = A\mathbf{x}_1 + A\mathbf{x}_2 = f_A(\mathbf{x}_1) + f_A(\mathbf{x}_2)$$

Second

What about $f_A(c\mathbf{x}_1)$?



Proof

If $A_{m \times n}$ is a matrix, f_A is a linear transformation

How?

First

$$f_A(\mathbf{x}_1 + \mathbf{x}_2) = A(\mathbf{x}_1 + \mathbf{x}_2) = A\mathbf{x}_1 + A\mathbf{x}_2 = f_A(\mathbf{x}_1) + f_A(\mathbf{x}_2)$$

Second

What about $f_A(c\mathbf{x}_1)$?



Outline

1 Linear Transformation

- Introduction
- Functions that can be defined using matrices
- Linear Functions
- **Kernel and Range**
- The Matrix of a Linear Transformation
- Going Back to Homogeneous Equations
- The Rank-Nullity Theorem

2 Derivative of Transformations

- Introduction
- Derivative of a Linear Transformation
- Derivative of a Quadratic Transformation

3 Linear Regression

- The Simplest Functions
- Splitting the Space
- Defining the Decision Surface
- Properties of the Hyperplane $w^T x + w_0$
- Augmenting the Vector
- Least Squared Error Procedure
 - The Geometry of a Two-Category Linearly-Separable Case
- The Error Idea
- The Final Error Equation

4 Principal Component Analysis

- Karhunen-Loeve Transform
- Projecting the Data
- Lagrange Multipliers
- The Process
- Example

5 Singular Value Decomposition

- Introduction
- Image Compression



We have

Definition (Actually related the null-space)

If $f : \mathbb{R}^n \rightarrow \mathbb{R}^m$ is linear, the kernel of f is defined by

$$\text{Ker}(f) = \{\mathbf{v} \in \mathbb{R}^n \mid f(\mathbf{v}) = 0\}$$

Definition

If $f : \mathbb{R}^n \rightarrow \mathbb{R}^m$ is linear, the range of f is defined by

$$\text{Range}(f) = \{\mathbf{y} \in \mathbb{R}^m \mid \mathbf{y} = f(\mathbf{x}) \text{ for some } \mathbf{x} \in \mathbb{R}^n\}$$



We have

Definition (Actually related the null-space)

If $f : \mathbb{R}^n \rightarrow \mathbb{R}^m$ is linear, the kernel of f is defined by

$$\text{Ker}(f) = \{\mathbf{v} \in \mathbb{R}^n \mid f(\mathbf{v}) = \mathbf{0}\}$$

Definition

If $f : \mathbb{R}^n \rightarrow \mathbb{R}^m$ is linear, the range of f is defined by

$$\text{Range}(f) = \{\mathbf{y} \in \mathbb{R}^m \mid \mathbf{y} = f(\mathbf{x}) \text{ for some } \mathbf{x} \in \mathbb{R}^n\}$$



We have also the following Spaces

Row Space

We have that the span of the row vectors of A form a subspace.

Column Space

We have that the span of the column vectors of A , also, form a subspace.



We have also the following Spaces

Row Space

We have that the span of the row vectors of A form a subspace.

Column Space

We have that the span of the column vectors of A , also, form a subspace.



From This

It can be shown that

$\text{Ker}(f)$ is a subspace of \mathbb{R}^n

Also

$\text{Range}(f)$ is a subspace of \mathbb{R}^m



From This

It can be shown that

$\text{Ker}(f)$ is a subspace of \mathbb{R}^n

Also

$\text{Range}(f)$ is a subspace of \mathbb{R}^m



Outline

1 Linear Transformation

- Introduction
- Functions that can be defined using matrices
- Linear Functions
- Kernel and Range
- **The Matrix of a Linear Transformation**
- Going Back to Homogeneous Equations
- The Rank-Nullity Theorem

2 Derivative of Transformations

- Introduction
- Derivative of a Linear Transformation
- Derivative of a Quadratic Transformation

3 Linear Regression

- The Simplest Functions
- Splitting the Space
- Defining the Decision Surface
- Properties of the Hyperplane $w^T x + w_0$
- Augmenting the Vector
- Least Squared Error Procedure
 - The Geometry of a Two-Category Linearly-Separable Case
- The Error Idea
- The Final Error Equation

4 Principal Component Analysis

- Karhunen-Loeve Transform
- Projecting the Data
- Lagrange Multipliers
- The Process
- Example

5 Singular Value Decomposition

- Introduction
- Image Compression

Assume the following

Let

$$e_1 = \begin{pmatrix} 1 \\ 0 \\ \vdots \\ 0 \end{pmatrix}_{n \times 1}, e_2 = \begin{pmatrix} 0 \\ 1 \\ \vdots \\ 0 \end{pmatrix}_{n \times 1}, \dots, e_n = \begin{pmatrix} 0 \\ 0 \\ \vdots \\ 1 \end{pmatrix}_{n \times 1}$$

Then any vector $x \in \mathbb{R}^n$

$$x = \begin{pmatrix} x_1 \\ x_2 \\ \vdots \\ x_n \end{pmatrix} = x_1 e_1 + x_2 e_2 + \dots + x_n e_n$$



Assume the following

Let

$$\mathbf{e}_1 = \begin{pmatrix} 1 \\ 0 \\ \vdots \\ 0 \end{pmatrix}_{n \times 1}, \mathbf{e}_3 = \begin{pmatrix} 0 \\ 1 \\ \vdots \\ 0 \end{pmatrix}_{n \times 1}, \dots, \mathbf{e}_n = \begin{pmatrix} 0 \\ 0 \\ \vdots \\ 1 \end{pmatrix}_{n \times 1}$$

Then any vector $\mathbf{x} \in \mathbb{R}^n$

$$\mathbf{x} = \begin{pmatrix} x_1 \\ x_2 \\ \vdots \\ x_n \end{pmatrix} = x_1 \mathbf{e}_1 + x_2 \mathbf{e}_2 + \dots + x_n \mathbf{e}_n$$



Then

Applying f

$$f(\mathbf{x}) = x_1 f(\mathbf{e}_1) + x_2 f(\mathbf{e}_2) + \dots + x_n f(\mathbf{e}_n)$$

A linear combination of elements

$$\{f(\mathbf{e}_1), f(\mathbf{e}_2), \dots, f(\mathbf{e}_n)\}$$

They are column vectors in \mathbb{R}^m

$$A = (f(\mathbf{e}_1) | f(\mathbf{e}_2) | \dots | f(\mathbf{e}_n))_{m \times n}$$



Then

Applying f

$$f(\mathbf{x}) = x_1 f(\mathbf{e}_1) + x_2 f(\mathbf{e}_2) + \dots + x_n f(\mathbf{e}_n)$$

A linear combination of elements

$$\{f(\mathbf{e}_1), f(\mathbf{e}_2), \dots, f(\mathbf{e}_n)\}$$

They are column vectors in

$$A = (f(\mathbf{e}_1) | f(\mathbf{e}_2) | \dots | f(\mathbf{e}_n))_{m \times n}$$



Then

Applying f

$$f(\mathbf{x}) = x_1 f(\mathbf{e}_1) + x_2 f(\mathbf{e}_2) + \dots + x_n f(\mathbf{e}_n)$$

A linear combination of elements

$$\{f(\mathbf{e}_1), f(\mathbf{e}_2), \dots, f(\mathbf{e}_n)\}$$

They are column vectors in \mathbb{R}^m

$$A = (f(\mathbf{e}_1) | f(\mathbf{e}_2) | \dots | f(\mathbf{e}_n))_{m \times n}$$



Thus, we have

Finally, we have

$$f(\mathbf{x}) = (f(\mathbf{e}_1) | f(\mathbf{e}_2) | \dots | f(\mathbf{e}_n)) \mathbf{x} = A\mathbf{x}$$

Definition

- The matrix A defined above for the function f is called the matrix of f in the standard basis.



Thus, we have

Finally, we have

$$f(\mathbf{x}) = (f(\mathbf{e}_1) | f(\mathbf{e}_2) | \dots | f(\mathbf{e}_n)) \mathbf{x} = A\mathbf{x}$$

Definition

- The matrix A defined above for the function f is called the matrix of f in the standard basis.



Outline

1 Linear Transformation

- Introduction
- Functions that can be defined using matrices
- Linear Functions
- Kernel and Range
- The Matrix of a Linear Transformation
- **Going Back to Homogeneous Equations**
- The Rank-Nullity Theorem

2 Derivative of Transformations

- Introduction
- Derivative of a Linear Transformation
- Derivative of a Quadratic Transformation

3 Linear Regression

- The Simplest Functions
- Splitting the Space
- Defining the Decision Surface
- Properties of the Hyperplane $w^T x + w_0$
- Augmenting the Vector
- Least Squared Error Procedure
 - The Geometry of a Two-Category Linearly-Separable Case
- The Error Idea
- The Final Error Equation

4 Principal Component Analysis

- Karhunen-Loeve Transform
- Projecting the Data
- Lagrange Multipliers
- The Process
- Example

5 Singular Value Decomposition

- Introduction
- Image Compression



Given an $m \times n$ matrix A

The set of all solutions to the homogeneous equation Ax

- It is a subspace V of \mathbb{R}^n .

$$Ax = 0$$

Remember how to prove the subspaces

$$x_1 + x_2 \in V \text{ and } cx \in V$$

- Do you remember?



Given an $m \times n$ matrix A

The set of all solutions to the homogeneous equation Ax

- It is a subspace V of \mathbb{R}^n .

$$Ax = 0$$

Remember how to prove the subspaces...

$$x_2 + x_2 \in V \text{ and } cx \in V$$

- Do you remember?



Then, we have

Definition

- This important subspace is called the null space of A , and is denoted $Null(A)$

It is also known as

$$x_H = \{x | Ax = 0\}$$



Then, we have

Definition

- This important subspace is called the null space of A , and is denoted $Null(A)$

It is also known as

$$\mathbf{x}_H = \{\mathbf{x} | A\mathbf{x} = 0\}$$



Knowing that $\text{Range}(f)$ and $\text{Ker}(f)$ are sub-spaces

Which ones they are?

Any Idea?

Range(f)

The column space of the matrix A .

Ker(f)

It is the null space of A .



Knowing that $\text{Range}(f)$ and $\text{Ker}(f)$ are sub-spaces

Which ones they are?

Any Idea?

$\text{Range}(f)$

The column space of the matrix A .

$\text{Ker}(f)$

It is the null space of A .



Knowing that $\text{Range}(f)$ and $\text{Ker}(f)$ are sub-spaces

Which ones they are?

Any Idea?

$\text{Range}(f)$

The column space of the matrix A .

$\text{Ker}(f)$

It is the null space of A .



Outline

1 Linear Transformation

- Introduction
- Functions that can be defined using matrices
- Linear Functions
- Kernel and Range
- The Matrix of a Linear Transformation
- Going Back to Homogeneous Equations
- **The Rank-Nullity Theorem**

2 Derivative of Transformations

- Introduction
- Derivative of a Linear Transformation
- Derivative of a Quadratic Transformation

3 Linear Regression

- The Simplest Functions
- Splitting the Space
- Defining the Decision Surface
- Properties of the Hyperplane $w^T x + w_0$
- Augmenting the Vector
- Least Squared Error Procedure
 - The Geometry of a Two-Category Linearly-Separable Case
- The Error Idea
- The Final Error Equation

4 Principal Component Analysis

- Karhunen-Loeve Transform
- Projecting the Data
- Lagrange Multipliers
- The Process
- Example

5 Singular Value Decomposition

- Introduction
- Image Compression



We have a nice theorem

Dimension Theorem

Let $f : \mathbb{R}^n \rightarrow \mathbb{R}^m$ be linear. Then

$$\dim(\text{domain}(f)) = \dim(\text{Range}(f)) + \dim(\text{Ker}(f))$$

Where

The dimension of V , written $\dim(V)$, is the number of elements in any basis of V .



We have a nice theorem

Dimension Theorem

Let $f : \mathbb{R}^n \rightarrow \mathbb{R}^m$ be linear. Then

$$\dim(\text{domain}(f)) = \dim(\text{Range}(f)) + \dim(\text{Ker}(f))$$

Where

The dimension of V , written $\dim(V)$, is the number of elements in any basis of V .



Rank and Nullity of a Matrix

Definition

- The rank of the matrix A is the dimension of the row space of A , and is denoted $R(A)$.

Example

- The rank of $I_{n \times n}$ is n .



Rank and Nullity of a Matrix

Definition

- The rank of the matrix A is the dimension of the row space of A , and is denoted $R(A)$.

Example

- The rank of $I_{n \times n}$ is n .



Then

Theorem

The rank of a matrix in Gauss-Jordan form is equal to the number of leading variables.

Proof

- In the G form of a matrix, every non-zero row has a leading 1, which is the only non-zero entry in its column.

Then

- No elementary row operation can zero out a leading 1, so these non-zero rows are linearly independent.



Then

Theorem

The rank of a matrix in Gauss-Jordan form is equal to the number of leading variables.

Proof

- In the G form of a matrix, every non-zero row has a leading 1, which is the only non-zero entry in its column.

Then

- No elementary row operation can zero out a leading 1, so these non-zero rows are linearly independent.



Then

Theorem

The rank of a matrix in Gauss-Jordan form is equal to the number of leading variables.

Proof

- In the G form of a matrix, every non-zero row has a leading 1, which is the only non-zero entry in its column.

Then

- No elementary row operation can zero out a leading 1, so these non-zero rows are linearly independent.



Thus

We have

- Since all the other rows are zero, the dimension of the row space of the Gauss-Jordan form is equal to the number of leading 1's.

Finally

- This is the same as the number of leading variables. Q.E.D.



Cinvestav

Thus

We have

- Since all the other rows are zero, the dimension of the row space of the Gauss-Jordan form is equal to the number of leading 1's.

Finally

- This is the same as the number of leading variables. Q.E.D.



Cinvestav

About the Nullity of the Matrix

Definition

- The nullity of the matrix A is the dimension of the null space of A , and is denoted by $\dim [N(A)]$.

Example

- The nullity of I is 0.



About the Nullity of the Matrix

Definition

- The nullity of the matrix A is the dimension of the null space of A , and is denoted by $\dim [N(A)]$.

Example

- The nullity of I is 0.



About the Nullity of the Matrix

Definition

- The nullity of the matrix A is the dimension of the null space of A , and is denoted by $\dim [N(A)]$.

Example

- The nullity of I is 0.



Number of Free Variables

Theorem

The nullity of a matrix in Gauss-Jordan form is equal to the number of free variables.

Proof

- Suppose A is $m \times n$, and that the Gauss-Jordan form has j leading variables and k free variables:

$$j + k = n$$



Proof

Then, when computing the solution to the homogeneous equation

- We solve for the first j (leading) variables in terms of the remaining k free variables:

$$s_1, s_2, s_3, \dots, s_k$$

Then

- Then the general solution to the homogeneous equation are:

$$s_1 \mathbf{v}_1 + s_2 \mathbf{v}_2 + s_3 \mathbf{v}_3 + \dots + s_k \mathbf{v}_k$$



Proof

Then, when computing the solution to the homogeneous equation

- We solve for the first j (leading) variables in terms of the remaining k free variables:

$$s_1, s_2, s_3, \dots, s_k$$

Then

- Then the general solution to the homogeneous equation are:

$$s_1\mathbf{v}_1 + s_2\mathbf{v}_2 + s_3\mathbf{v}_3 + \dots + s_k\mathbf{v}_k$$



Proof

Then, when computing the solution to the homogeneous equation

- We solve for the first j (leading) variables in terms of the remaining k free variables:

$$s_1, s_2, s_3, \dots, s_k$$

Then

- Then the general solution to the homogeneous equation are:

$$s_1\mathbf{v}_1 + s_2\mathbf{v}_2 + s_3\mathbf{v}_3 + \dots + s_k\mathbf{v}_k$$



Where

The vectors are the Canonical Ones

- Here, a trick!!!

Meaning in v_j , we have 1, after many 0

- It appears at position $j + 1$, with zeros afterwards, and so on.

Therefore the vectors are linearly independent

$$v_1, v_2, v_3, \dots, v_k$$



Where

The vectors are the Canonical Ones

- Here, a trick!!!

Meaning in v_1 , we have 1, after many 0

- It appears at position $j + 1$, with zeros afterwards, and so on.

Therefore, the vectors are linearly independent:

$$v_1, v_2, v_3, \dots, v_k$$



Where

The vectors are the Canonical Ones

- Here, a trick!!!

Meaning in v_1 , we have 1, after many 0

- It appears at position $j + 1$, with zeros afterwards, and so on.

Therefore the vectors are linearly independents

$$v_1, v_2, v_3, \dots, v_k$$



Therefore

They are a basis for the null space of A

And there are k of them, the same as the number of free variables.



Cinvestav

Definition

The matrix B is said to be row equivalent to A ($B \sim A$) if B can be obtained from A by a finite sequence of elementary row operations.

Equivalent Form:

$B \sim A \Leftrightarrow$ There exist elementary matrices such that

$$B = E_k E_{k-1} E_{k-2} \cdots E_1 A$$

If we write $C = E_k E_{k-1} E_{k-2} \cdots E_1$

B is row equivalent to A if $B = CA$ with C invertible.



Now

Definition

The matrix B is said to be row equivalent to A ($B \sim A$) if B can be obtained from A by a finite sequence of elementary row operations.

In matrix terms

$B \sim A \Leftrightarrow$ There exist elementary matrices such that

$$B = E_k E_{k-1} E_{k-1} \cdots E_1 A$$

If we write $E = E_k \cdots E_1$

B is row equivalent to A if $B = CA$ with C invertible.



Now

Definition

The matrix B is said to be row equivalent to A ($B \sim A$) if B can be obtained from A by a finite sequence of elementary row operations.

In matrix terms

$B \sim A \Leftrightarrow$ There exist elementary matrices such that

$$B = E_k E_{k-1} E_{k-1} \cdots E_1 A$$

If we write $C = E_k E_{k-1} E_{k-1} \cdots E_1$

B is row equivalent to A if $B = CA$ with C invertible.



Then, we have

Theorem

If $B \sim A$, then $\text{Null}(B) = \text{Null}(A)$.

Theorem

If $B \sim A$, then the row space of B is identical to that of A .

Summarizing

Row operations change neither the row space nor the null space of A .



Then, we have

Theorem

If $B \sim A$, then $\text{Null}(B) = \text{Null}(A)$.

Theorem

If $B \sim A$, then the row space of B is identical to that of A .

Summary

Row operations change neither the row space nor the null space of A .



Then, we have

Theorem

If $B \sim A$, then $\text{Null}(B) = \text{Null}(A)$.

Theorem

If $B \sim A$, then the row space of B is identical to that of A .

Summarizing

Row operations change neither the row space nor the null space of A .



Corollaries

Corollary 1

- If R is the Gauss-Jordan form of A , then R has the same null space and row space as A .

Corollary 2

- If $B \sim A$, then $R(B) = R(A)$, and $N(B) = N(A)$.



Corollaries

Corollary 1

- If R is the Gauss-Jordan form of A , then R has the same null space and row space as A .

Corollary 2

- If $B \sim A$, then $R(B) = R(A)$, and $N(B) = N(A)$.



Then

Theorem

- The number of linearly independent rows of the matrix A is equal to the number of linearly independent columns of A .

Thus:

- In particular, the rank of A is also equal to the number of linearly independent columns, and hence to the dimension of the column space of A .

Therefore:

- The number of linearly independent columns of A is then just the number of leading entries in the Gauss-Jordan form of A which is, as we know, the same as the rank of A .

Then

Theorem

- The number of linearly independent rows of the matrix A is equal to the number of linearly independent columns of A .

Thus

- In particular, the rank of A is also equal to the number of linearly independent columns, and hence to the dimension of the column space of A

Therefore

- The number of linearly independent columns of A is then just the number of leading entries in the Gauss-Jordan form of A which is, as we know, the same as the rank of A .

Then

Theorem

- The number of linearly independent rows of the matrix A is equal to the number of linearly independent columns of A .

Thus

- In particular, the rank of A is also equal to the number of linearly independent columns, and hence to the dimension of the column space of A

Therefore

- The number of linearly independent columns of A is then just the number of leading entries in the Gauss-Jordan form of A which is, as we know, the same as the rank of A .

Proof of the theorem (Dimension Theorem)

First

- The rank of A is the same as the rank of the Gauss-Jordan form of A which is equal to the number of leading entries in the Gauss-Jordan form.

Additionally

- The dimension of the null space is equal to the number of free variables in the reduced echelon (Gauss-Jordan) form of A .

Then

We know further that the number of free variables plus the number of leading entries is exactly the number of columns.



Proof of the theorem (Dimension Theorem)

First

- The rank of A is the same as the rank of the Gauss-Jordan form of A which is equal to the number of leading entries in the Gauss-Jordan form.

Additionally

- The dimension of the null space is equal to the number of free variables in the reduced echelon (Gauss-Jordan) form of A .

Proof

We know further that the number of free variables plus the number of leading entries is exactly the number of columns.



Proof of the theorem (Dimension Theorem)

First

- The rank of A is the same as the rank of the Gauss-Jordan form of A which is equal to the number of leading entries in the Gauss-Jordan form.

Additionally

- The dimension of the null space is equal to the number of free variables in the reduced echelon (Gauss-Jordan) form of A .

Then

We know further that the number of free variables plus the number of leading entries is exactly the number of columns.



Finally

We have

$$\dim(\text{domain}(f)) = \dim(\text{Range}(f)) + \dim(\text{Ker}(f))$$



Cinvestav

Outline

1 Linear Transformation

- Introduction
- Functions that can be defined using matrices
- Linear Functions
- Kernel and Range
- The Matrix of a Linear Transformation
- Going Back to Homogeneous Equations
- The Rank-Nullity Theorem

2 Derivative of Transformations

- **Introduction**
- Derivative of a Linear Transformation
- Derivative of a Quadratic Transformation

3 Linear Regression

- The Simplest Functions
- Splitting the Space
- Defining the Decision Surface
- Properties of the Hyperplane $w^T x + w_0$
- Augmenting the Vector
- Least Squared Error Procedure
 - The Geometry of a Two-Category Linearly-Separable Case
- The Error Idea
- The Final Error Equation

4 Principal Component Analysis

- Karhunen-Loeve Transform
- Projecting the Data
- Lagrange Multipliers
- The Process
- Example

5 Singular Value Decomposition

- Introduction
- Image Compression



As we know

Many Times

We want to obtain a maximum or a minimum of a cost function expressed in terms of matrices....

We need then to define matrix derivatives

Thus, this discussion is useful in Machine Learning.



Cinvestav

As we know

Many Times

We want to obtain a maximum or a minimum of a cost function expressed in terms of matrices....

We need then to define matrix derivatives

Thus, this discussion is useful in Machine Learning.



Cinvestav

Basic Definition

Let $\psi(\mathbf{x}) = \mathbf{y}$

Where \mathbf{y} is an m -element vector, and \mathbf{x} is an n -element vector

Then, we define the derivative with respect to a vector

$$\frac{\partial \mathbf{y}}{\partial \mathbf{x}} = \begin{bmatrix} \frac{\partial y_1}{\partial x_1} & \frac{\partial y_1}{\partial x_2} & \dots & \frac{\partial y_1}{\partial x_n} \\ \frac{\partial y_2}{\partial x_1} & \frac{\partial y_2}{\partial x_2} & \dots & \frac{\partial y_2}{\partial x_n} \\ \vdots & \vdots & \ddots & \vdots \\ \frac{\partial y_m}{\partial x_1} & \frac{\partial y_m}{\partial x_2} & \dots & \frac{\partial y_m}{\partial x_n} \end{bmatrix}$$



Basic Definition

Let $\psi(\mathbf{x}) = \mathbf{y}$

Where \mathbf{y} is an m -element vector, and \mathbf{x} is an n -element vector

Then, we define the derivative with respect to a vector

$$\frac{\partial \mathbf{y}}{\partial \mathbf{x}} = \begin{bmatrix} \frac{\partial y_1}{\partial x_1} & \frac{\partial y_1}{\partial x_2} & \dots & \frac{\partial y_1}{\partial x_n} \\ \frac{\partial y_2}{\partial x_1} & \frac{\partial y_2}{\partial x_2} & \dots & \frac{\partial y_2}{\partial x_n} \\ \vdots & \vdots & \ddots & \vdots \\ \frac{\partial y_m}{\partial x_1} & \frac{\partial y_m}{\partial x_2} & \dots & \frac{\partial y_m}{\partial x_n} \end{bmatrix}$$



What is this

The Matrix denotes the $m \times n$ matrix of first order partial derivatives

- Such a matrix is called the Jacobian matrix of the transformation $\psi(\mathbf{x})$.

When we can get our first ideas on derivatives

- For Linear Transformations.



What is this

The Matrix denotes the $m \times n$ matrix of first order partial derivatives

- Such a matrix is called the Jacobian matrix of the transformation $\psi(\mathbf{x})$.

Then, we can get our first ideas on derivatives

- For Linear Transformations.



Outline

1 Linear Transformation

- Introduction
- Functions that can be defined using matrices
- Linear Functions
- Kernel and Range
- The Matrix of a Linear Transformation
- Going Back to Homogeneous Equations
- The Rank-Nullity Theorem

2 Derivative of Transformations

- Introduction
- **Derivative of a Linear Transformation**
- Derivative of a Quadratic Transformation

3 Linear Regression

- The Simplest Functions
- Splitting the Space
- Defining the Decision Surface
- Properties of the Hyperplane $w^T x + w_0$
- Augmenting the Vector
- Least Squared Error Procedure
 - The Geometry of a Two-Category Linearly-Separable Case
- The Error Idea
- The Final Error Equation

4 Principal Component Analysis

- Karhunen-Loeve Transform
- Projecting the Data
- Lagrange Multipliers
- The Process
- Example

5 Singular Value Decomposition

- Introduction
- Image Compression



Derivative of $\mathbf{y} = A\mathbf{x}$

Theorem

- Let $\mathbf{y} = A\mathbf{x}$ where \mathbf{y} is a $m \times 1$, \mathbf{x} is a $n \times 1$, A is a $m \times n$ and A does not depend on \mathbf{x} , then

$$\frac{\partial \mathbf{y}}{\partial \mathbf{x}} = A$$



Proof

Each i^{th} element of \mathbf{y} is given by

$$y_i = \sum_{k=1}^N a_{ik}x_k$$

We have that

$$\frac{\partial y_i}{\partial x_j} = a_{ij}$$

for all $i = 1, 2, \dots, m$ and $j = 1, 2, \dots, n$

Hence

$$\frac{\partial \mathbf{y}}{\partial \mathbf{x}} = \mathbf{A}$$

Proof

Each i^{th} element of \mathbf{y} is given by

$$y_i = \sum_{k=1}^N a_{ik}x_k$$

We have that

$$\frac{\partial y_i}{\partial x_j} = a_{ij}$$

for all $i = 1, 2, \dots, m$ and $j = 1, 2, \dots, n$

Hence

$$\frac{\partial \mathbf{y}}{\partial \mathbf{x}} = \mathbf{A}$$

Proof

Each i^{th} element of \mathbf{y} is given by

$$y_i = \sum_{k=1}^N a_{ik}x_k$$

We have that

$$\frac{\partial y_i}{\partial x_j} = a_{ij}$$

for all $i = 1, 2, \dots, m$ and $j = 1, 2, \dots, n$

Hence

$$\frac{\partial \mathbf{y}}{\partial \mathbf{x}} = A$$

Outline

1 Linear Transformation

- Introduction
- Functions that can be defined using matrices
- Linear Functions
- Kernel and Range
- The Matrix of a Linear Transformation
- Going Back to Homogeneous Equations
- The Rank-Nullity Theorem

2 Derivative of Transformations

- Introduction
- Derivative of a Linear Transformation
- **Derivative of a Quadratic Transformation**

3 Linear Regression

- The Simplest Functions
- Splitting the Space
- Defining the Decision Surface
- Properties of the Hyperplane $w^T x + w_0$
- Augmenting the Vector
- Least Squared Error Procedure
 - The Geometry of a Two-Category Linearly-Separable Case
- The Error Idea
- The Final Error Equation

4 Principal Component Analysis

- Karhunen-Loeve Transform
- Projecting the Data
- Lagrange Multipliers
- The Process
- Example

5 Singular Value Decomposition

- Introduction
- Image Compression



Derivative of $\mathbf{y}^T A \mathbf{x}$

Theorem

- Let the scalar α be defined by

$$\alpha = \mathbf{y}^T A \mathbf{x}$$

where

\mathbf{y} is a $m \times 1$, \mathbf{x} is a $n \times 1$, A is a $m \times n$ and A does not depend on \mathbf{x} and \mathbf{y} , then

$$\frac{\partial \alpha}{\partial \mathbf{x}} = \mathbf{y}^T A \text{ and } \frac{\partial \alpha}{\partial \mathbf{y}} = \mathbf{x}^T A^T$$



Derivative of $\mathbf{y}^T A \mathbf{x}$

Theorem

- Let the scalar α be defined by

$$\alpha = \mathbf{y}^T A \mathbf{x}$$

where

\mathbf{y} is a $m \times 1$, \mathbf{x} is a $n \times 1$, A is a $m \times n$ and A does not depend on \mathbf{x} and \mathbf{y} , then

$$\frac{\partial \alpha}{\partial \mathbf{x}} = \mathbf{y}^T A \text{ and } \frac{\partial \alpha}{\partial \mathbf{y}} = \mathbf{x}^T A^T$$



Proof

Define

$$\mathbf{w}^T = \mathbf{y}^T A$$

Note

$$\alpha = \mathbf{w}^T \mathbf{x}$$

By the previous theorem

$$\frac{\partial \alpha}{\partial \mathbf{x}} = \mathbf{w}^T = \mathbf{y}^T A$$

In a similar way, you can prove the other statement.



Proof

Define

$$\mathbf{w}^T = \mathbf{y}^T A$$

Note

$$\alpha = \mathbf{w}^T \mathbf{x}$$

By the previous theorem

$$\frac{\partial \alpha}{\partial \mathbf{x}} = \mathbf{w}^T = \mathbf{y}^T A$$

In a similar way, you can prove the other statement.



Proof

Define

$$\mathbf{w}^T = \mathbf{y}^T A$$

Note

$$\alpha = \mathbf{w}^T \mathbf{x}$$

By the previous theorem

$$\frac{\partial \alpha}{\partial \mathbf{x}} = \mathbf{w}^T = \mathbf{y}^T A$$

In a similar way, you can prove the other statement.



Outline

1 Linear Transformation

- Introduction
- Functions that can be defined using matrices
- Linear Functions
- Kernel and Range
- The Matrix of a Linear Transformation
- Going Back to Homogeneous Equations
- The Rank-Nullity Theorem

2 Derivative of Transformations

- Introduction
- Derivative of a Linear Transformation
- Derivative of a Quadratic Transformation

3 Linear Regression

- **The Simplest Functions**
- Splitting the Space
- Defining the Decision Surface
- Properties of the Hyperplane $w^T x + w_0$
- Augmenting the Vector
- Least Squared Error Procedure
 - The Geometry of a Two-Category Linearly-Separable Case
- The Error Idea
- The Final Error Equation

4 Principal Component Analysis

- Karhunen-Loeve Transform
- Projecting the Data
- Lagrange Multipliers
- The Process
- Example

5 Singular Value Decomposition

- Introduction
- Image Compression



What is it?

First than anything, we have a parametric model!!!

Here, we have an hyperplane as a model:

$$g(\mathbf{x}) = \mathbf{w}^T \mathbf{x} + w_0 \quad (1)$$

Note: $\mathbf{w}^T \mathbf{x}$ is also know as dot product

In the case of 2D:

We have:

$$g(\mathbf{x}) = (w_1, w_2) \begin{pmatrix} x_1 \\ x_2 \end{pmatrix} + w_0 = w_1 x_1 + w_2 x_2 + w_0 \quad (2)$$



What is it?

First than anything, we have a parametric model!!!

Here, we have an hyperplane as a model:

$$g(\mathbf{x}) = \mathbf{w}^T \mathbf{x} + w_0 \quad (1)$$

Note: $\mathbf{w}^T \mathbf{x}$ is also know as dot product

In the case of \mathbb{R}^2

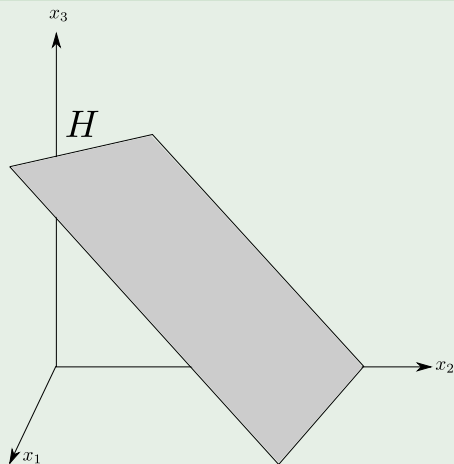
We have:

$$g(\mathbf{x}) = (w_1, w_2) \begin{pmatrix} x_1 \\ x_2 \end{pmatrix} + w_0 = w_1 x_1 + w_2 x_2 + w_0 \quad (2)$$



Example

Hyperplane in \mathbb{R}^3



Outline

1 Linear Transformation

- Introduction
- Functions that can be defined using matrices
- Linear Functions
- Kernel and Range
- The Matrix of a Linear Transformation
- Going Back to Homogeneous Equations
- The Rank-Nullity Theorem

2 Derivative of Transformations

- Introduction
- Derivative of a Linear Transformation
- Derivative of a Quadratic Transformation

3 Linear Regression

- The Simplest Functions
- **Splitting the Space**
- Defining the Decision Surface
- Properties of the Hyperplane $w^T x + w_0$
- Augmenting the Vector
- Least Squared Error Procedure
 - The Geometry of a Two-Category Linearly-Separable Case
- The Error Idea
- The Final Error Equation

4 Principal Component Analysis

- Karhunen-Loeve Transform
- Projecting the Data
- Lagrange Multipliers
- The Process
- Example

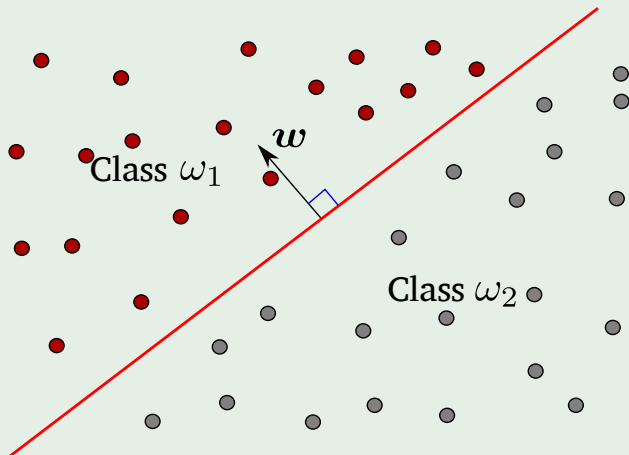
5 Singular Value Decomposition

- Introduction
- Image Compression



Splitting The Space \mathbb{R}^2

Using a simple straight line (Hyperplane)



Splitting the Space?

For example, assume the following vector w and constant w_0

$$w = (-1, 2)^T \text{ and } w_0 = 0$$

Hyperplane

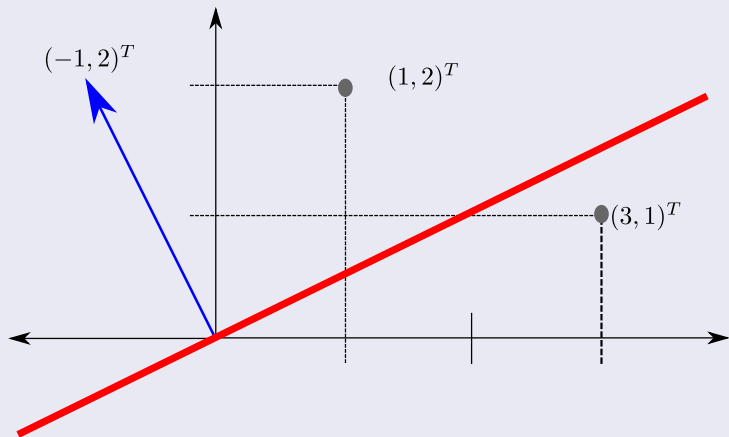


Splitting the Space?

For example, assume the following vector w and constant w_0

$$w = (-1, 2)^T \text{ and } w_0 = 0$$

Hyperplane



Then, we have

The following results

$$g\left(\begin{pmatrix} 1 \\ 2 \end{pmatrix}\right) = (-1, 2) \begin{pmatrix} 1 \\ 2 \end{pmatrix} = -1 \times 1 + 2 \times 2 = 3$$

$$g\left(\begin{pmatrix} 3 \\ 1 \end{pmatrix}\right) = (-1, 2) \begin{pmatrix} 3 \\ 1 \end{pmatrix} = -1 \times 3 + 2 \times 1 = -1$$

YES!!! We have a positive side and a negative side!!!



Outline

1 Linear Transformation

- Introduction
- Functions that can be defined using matrices
- Linear Functions
- Kernel and Range
- The Matrix of a Linear Transformation
- Going Back to Homogeneous Equations
- The Rank-Nullity Theorem

2 Derivative of Transformations

- Introduction
- Derivative of a Linear Transformation
- Derivative of a Quadratic Transformation

3 Linear Regression

- The Simplest Functions
- Splitting the Space
- **Defining the Decision Surface**
- Properties of the Hyperplane $w^T x + w_0$
- Augmenting the Vector
- Least Squared Error Procedure
 - The Geometry of a Two-Category Linearly-Separable Case
- The Error Idea
- The Final Error Equation

4 Principal Component Analysis

- Karhunen-Loeve Transform
- Projecting the Data
- Lagrange Multipliers
- The Process
- Example

5 Singular Value Decomposition

- Introduction
- Image Compression



The Decision Surface

The equation $g(x) = 0$ defines a decision surface

Separating the elements in classes, ω_1 and ω_2 .

When $g(x)$ is linear the decision surface is an hyperplane

Now assume x_1 and x_2 are both on the decision surface

$$w^T x_1 + w_0 = 0$$

$$w^T x_2 + w_0 = 0$$

This

$$w^T x_1 + w_0 = w^T x_2 + w_0$$

(3)

The Decision Surface

The equation $g(x) = 0$ defines a decision surface

Separating the elements in classes, ω_1 and ω_2 .

When $g(x)$ is linear the decision surface is an hyperplane

Now assume x_1 and x_2 are both on the decision surface

$$w^T x_1 + w_0 = 0$$

$$w^T x_2 + w_0 = 0$$

This

$$w^T x_1 + w_0 = w^T x_2 + w_0$$

(3)

The Decision Surface

The equation $g(x) = 0$ defines a decision surface

Separating the elements in classes, ω_1 and ω_2 .

When $g(x)$ is linear the decision surface is an hyperplane

Now assume x_1 and x_2 are both on the decision surface

$$\mathbf{w}^T \mathbf{x}_1 + w_0 = 0$$

$$\mathbf{w}^T \mathbf{x}_2 + w_0 = 0$$

Thus

$$\mathbf{w}^T \mathbf{x}_1 + w_0 = \mathbf{w}^T \mathbf{x}_2 + w_0 \quad (3)$$

Defining a Decision Surface

Then

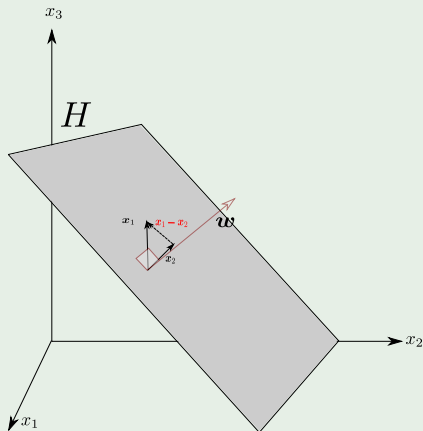
$$\mathbf{w}^T (\mathbf{x}_1 - \mathbf{x}_2) = 0 \quad (4)$$



Therefore

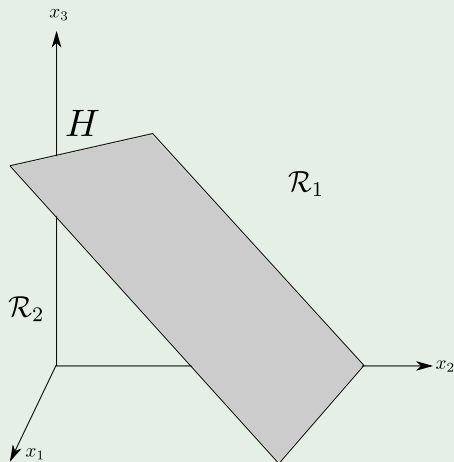
$x_1 - x_2$ lives in the hyperplane i.e. it is perpendicular to w^T

- Remark: any vector in the hyperplane is a linear combination of elements in a basis
- **Therefore any vector in the plane is perpendicular to w^T**



Therefore

The space is split in two regions (Example in \mathbb{R}^3) by the hyperplane H



Outline

1 Linear Transformation

- Introduction
- Functions that can be defined using matrices
- Linear Functions
- Kernel and Range
- The Matrix of a Linear Transformation
- Going Back to Homogeneous Equations
- The Rank-Nullity Theorem

2 Derivative of Transformations

- Introduction
- Derivative of a Linear Transformation
- Derivative of a Quadratic Transformation

3 Linear Regression

- The Simplest Functions
- Splitting the Space
- Defining the Decision Surface
- **Properties of the Hyperplane $w^T x + w_0$**
- Augmenting the Vector
- Least Squared Error Procedure
 - The Geometry of a Two-Category Linearly-Separable Case
- The Error Idea
- The Final Error Equation

4 Principal Component Analysis

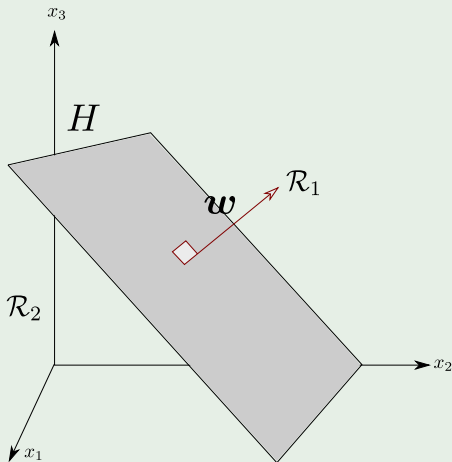
- Karhunen-Loeve Transform
- Projecting the Data
- Lagrange Multipliers
- The Process
- Example

5 Singular Value Decomposition

- Introduction
- Image Compression

Some Properties of the Hyperplane

Given that $g(\mathbf{x}) > 0$ if $\mathbf{x} \in \mathcal{R}_1$



It is more

We can say the following

- Any $x \in \mathcal{R}_1$ is on the positive side of H .
- Any $x \in \mathcal{R}_2$ is on the negative side of H .

It is more

We can say the following

- Any $x \in \mathcal{R}_1$ is on the positive side of H .
- Any $x \in \mathcal{R}_2$ is on the negative side of H .

In addition, $y(x)$ can give us a way to obtain the distance from x to the hyperplane H .

First, we express any x as follows

$$x = x_p + r \frac{w}{\|w\|}$$

It is more

We can say the following

- Any $x \in \mathcal{R}_1$ is on the positive side of H .
- Any $x \in \mathcal{R}_2$ is on the negative side of H .

In addition, $g(x)$ can give us a way to obtain the distance from x to the hyperplane H

First, we express any x as follows

$$\mathbf{x} = \mathbf{x}_p + r \frac{\mathbf{w}}{\|\mathbf{w}\|}$$

where

- \mathbf{x}_p is the normal projection of x onto H .
- r is the desired distance
 - ▶ Positive, if x is in the positive side
 - ▶ Negative, if x is in the negative side

It is more

We can say the following

- Any $x \in \mathcal{R}_1$ is on the positive side of H .
- Any $x \in \mathcal{R}_2$ is on the negative side of H .

In addition, $g(x)$ can give us a way to obtain the distance from x to the hyperplane H

First, we express any x as follows

$$\mathbf{x} = \mathbf{x}_p + r \frac{\mathbf{w}}{\|\mathbf{w}\|}$$

Where

- \mathbf{x}_p is the normal projection of x onto H .
- r is the desired distance
 - ▶ Positive, if x is in the positive side
 - ▶ Negative, if x is in the negative side

It is more

We can say the following

- Any $x \in \mathcal{R}_1$ is on the positive side of H .
- Any $x \in \mathcal{R}_2$ is on the negative side of H .

In addition, $g(x)$ can give us a way to obtain the distance from x to the hyperplane H

First, we express any x as follows

$$x = x_p + r \frac{w}{\|w\|}$$

Where

- x_p is the normal projection of x onto H .
- r is the desired distance

► Positive, if x is in the positive side
► Negative, if x is in the negative side

It is more

We can say the following

- Any $x \in \mathcal{R}_1$ is on the positive side of H .
- Any $x \in \mathcal{R}_2$ is on the negative side of H .

In addition, $g(x)$ can give us a way to obtain the distance from x to the hyperplane H

First, we express any x as follows

$$x = x_p + r \frac{w}{\|w\|}$$

Where

- x_p is the normal projection of x onto H .
- r is the desired distance
 - ▶ Positive, if x is in the positive side
 - ▶ Negative, if x is in the negative side

It is more

We can say the following

- Any $x \in \mathcal{R}_1$ is on the positive side of H .
- Any $x \in \mathcal{R}_2$ is on the negative side of H .

In addition, $g(x)$ can give us a way to obtain the distance from x to the hyperplane H

First, we express any x as follows

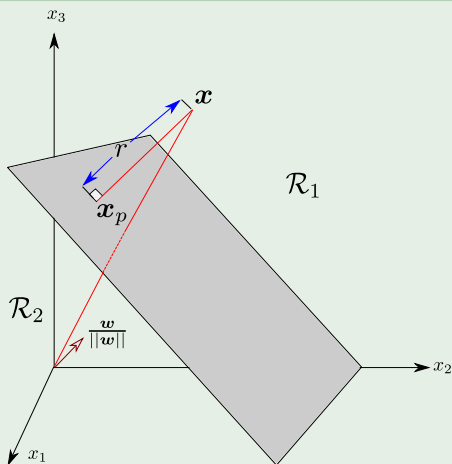
$$x = x_p + r \frac{w}{\|w\|}$$

Where

- x_p is the normal projection of x onto H .
- r is the desired distance
 - ▶ Positive, if x is in the positive side
 - ▶ Negative, if x is in the negative side

We have something like this

We have then



Now

Since $g(\mathbf{x}_p) = 0$

We have that

$$\begin{aligned}g(\mathbf{x}) &= g\left(\mathbf{x}_p + r \frac{\mathbf{w}}{\|\mathbf{w}\|}\right) \\&= \mathbf{w}^T \left(\mathbf{x}_p + r \frac{\mathbf{w}}{\|\mathbf{w}\|}\right) + w_0 \\&= \mathbf{w}^T \mathbf{x}_p + w_0 + r \frac{\mathbf{w}^T \mathbf{w}}{\|\mathbf{w}\|} \\&= g(\mathbf{x}_p) + r \frac{\|\mathbf{w}\|^2}{\|\mathbf{w}\|} \\&= r \|\mathbf{w}\|\end{aligned}$$

Now

Since $g(\mathbf{x}_p) = 0$

We have that

$$\begin{aligned}g(\mathbf{x}) &= g\left(\mathbf{x}_p + r \frac{\mathbf{w}}{\|\mathbf{w}\|}\right) \\&= \mathbf{w}^T \left(\mathbf{x}_p + r \frac{\mathbf{w}}{\|\mathbf{w}\|}\right) + w_0 \\&= \mathbf{w}^T \mathbf{x}_p + w_0 + r \frac{\mathbf{w}^T \mathbf{w}}{\|\mathbf{w}\|} \\&= g(\mathbf{x}_p) + r \frac{\|\mathbf{w}\|^2}{\|\mathbf{w}\|} \\&= r \|\mathbf{w}\|\end{aligned}$$

Then, we have

$$r = \frac{g(\mathbf{x})}{\|\mathbf{w}\|} \quad (5)$$

Now

Since $g(\mathbf{x}_p) = 0$

We have that

$$\begin{aligned}g(\mathbf{x}) &= g\left(\mathbf{x}_p + r \frac{\mathbf{w}}{\|\mathbf{w}\|}\right) \\&= \mathbf{w}^T \left(\mathbf{x}_p + r \frac{\mathbf{w}}{\|\mathbf{w}\|}\right) + w_0 \\&= \mathbf{w}^T \mathbf{x}_p + w_0 + r \frac{\mathbf{w}^T \mathbf{w}}{\|\mathbf{w}\|} \\&= g(\mathbf{x}_p) + r \frac{\|\mathbf{w}\|^2}{\|\mathbf{w}\|} \\&= r \|\mathbf{w}\|\end{aligned}$$

Then, we have

$$r = \frac{g(\mathbf{x})}{\|\mathbf{w}\|} \quad (5)$$

Now

Since $g(\mathbf{x}_p) = 0$

We have that

$$\begin{aligned}g(\mathbf{x}) &= g\left(\mathbf{x}_p + r \frac{\mathbf{w}}{\|\mathbf{w}\|}\right) \\&= \mathbf{w}^T \left(\mathbf{x}_p + r \frac{\mathbf{w}}{\|\mathbf{w}\|}\right) + w_0 \\&= \mathbf{w}^T \mathbf{x}_p + w_0 + r \frac{\mathbf{w}^T \mathbf{w}}{\|\mathbf{w}\|} \\&= g(\mathbf{x}_p) + r \frac{\|\mathbf{w}\|^2}{\|\mathbf{w}\|} \\&= r \|\mathbf{w}\|\end{aligned}$$

Then, we have

$$r = \frac{g(\mathbf{x})}{\|\mathbf{w}\|} \quad (5)$$

Now

Since $g(\mathbf{x}_p) = 0$

We have that

$$\begin{aligned}g(\mathbf{x}) &= g\left(\mathbf{x}_p + r \frac{\mathbf{w}}{\|\mathbf{w}\|}\right) \\&= \mathbf{w}^T \left(\mathbf{x}_p + r \frac{\mathbf{w}}{\|\mathbf{w}\|}\right) + w_0 \\&= \mathbf{w}^T \mathbf{x}_p + w_0 + r \frac{\mathbf{w}^T \mathbf{w}}{\|\mathbf{w}\|} \\&= g(\mathbf{x}_p) + r \frac{\|\mathbf{w}\|^2}{\|\mathbf{w}\|} \\&= r \|\mathbf{w}\|\end{aligned}$$

Then, we have

$$r = \frac{g(\mathbf{x})}{\|\mathbf{w}\|} \quad (5)$$

Now

Since $g(\mathbf{x}_p) = 0$

We have that

$$\begin{aligned}g(\mathbf{x}) &= g\left(\mathbf{x}_p + r \frac{\mathbf{w}}{\|\mathbf{w}\|}\right) \\&= \mathbf{w}^T \left(\mathbf{x}_p + r \frac{\mathbf{w}}{\|\mathbf{w}\|}\right) + w_0 \\&= \mathbf{w}^T \mathbf{x}_p + w_0 + r \frac{\mathbf{w}^T \mathbf{w}}{\|\mathbf{w}\|} \\&= g(\mathbf{x}_p) + r \frac{\|\mathbf{w}\|^2}{\|\mathbf{w}\|} \\&= r \|\mathbf{w}\|\end{aligned}$$

Then, we have

$$r = \frac{g(\mathbf{x})}{\|\mathbf{w}\|} \quad (5)$$

In particular

The distance from the origin to H

$$r = \frac{g(\mathbf{0})}{\|\mathbf{w}\|} = \frac{\mathbf{w}^T(\mathbf{0}) + w_0}{\|\mathbf{w}\|} = \frac{w_0}{\|\mathbf{w}\|} \quad (6)$$



In particular

The distance from the origin to H

$$r = \frac{g(\mathbf{0})}{\|\mathbf{w}\|} = \frac{\mathbf{w}^T(\mathbf{0}) + w_0}{\|\mathbf{w}\|} = \frac{w_0}{\|\mathbf{w}\|} \quad (6)$$

Remarks

- If $w_0 > 0$, the origin is on the positive side of H .
- If $w_0 < 0$, the origin is on the negative side of H .
- If $w_0 = 0$, the hyperplane has the homogeneous form $\mathbf{w}^T \mathbf{x}$ and hyperplane passes through the origin.



In particular

The distance from the origin to H

$$r = \frac{g(\mathbf{0})}{\|\mathbf{w}\|} = \frac{\mathbf{w}^T(\mathbf{0}) + w_0}{\|\mathbf{w}\|} = \frac{w_0}{\|\mathbf{w}\|} \quad (6)$$

Remarks

- If $w_0 > 0$, the origin is on the positive side of H .
- If $w_0 < 0$, the origin is on the negative side of H .
- If $w_0 = 0$, the hyperplane has the homogeneous form $\mathbf{w}^T \mathbf{x}$ and hyperplane passes through the origin.



In particular

The distance from the origin to H

$$r = \frac{g(\mathbf{0})}{\|\mathbf{w}\|} = \frac{\mathbf{w}^T(\mathbf{0}) + w_0}{\|\mathbf{w}\|} = \frac{w_0}{\|\mathbf{w}\|} \quad (6)$$

Remarks

- If $w_0 > 0$, the origin is on the positive side of H .
- If $w_0 < 0$, the origin is on the negative side of H .
- If $w_0 = 0$, the hyperplane has the homogeneous form $\mathbf{w}^T \mathbf{x}$ and hyperplane passes through the origin.



Outline

1 Linear Transformation

- Introduction
- Functions that can be defined using matrices
- Linear Functions
- Kernel and Range
- The Matrix of a Linear Transformation
- Going Back to Homogeneous Equations
- The Rank-Nullity Theorem

2 Derivative of Transformations

- Introduction
- Derivative of a Linear Transformation
- Derivative of a Quadratic Transformation

3 Linear Regression

- The Simplest Functions
- Splitting the Space
- Defining the Decision Surface
- Properties of the Hyperplane $w^T x + w_0$
- **Augmenting the Vector**
- Least Squared Error Procedure
 - The Geometry of a Two-Category Linearly-Separable Case
- The Error Idea
- The Final Error Equation

4 Principal Component Analysis

- Karhunen-Loeve Transform
- Projecting the Data
- Lagrange Multipliers
- The Process
- Example

5 Singular Value Decomposition

- Introduction
- Image Compression



We want to solve the independence of w_0

We would like w_0 as part of the dot product by making $x_0 = 1$

$$g(\mathbf{x}) = w_0 \times 1 + \sum_{i=1}^d w_i x_i = w_0 \times x_0 + \sum_{i=1}^d w_i x_i = \sum_{i=0}^d w_i x_i \quad (7)$$

We want to solve the independence of w_0

We would like w_0 as part of the dot product by making $x_0 = 1$

$$g(\mathbf{x}) = w_0 \times 1 + \sum_{i=1}^d w_i x_i = w_0 \times x_0 + \sum_{i=1}^d w_i x_i = \sum_{i=0}^d w_i x_i \quad (7)$$

By making

$$\mathbf{x}_{aug} = \begin{pmatrix} 1 \\ x_1 \\ \vdots \\ x_d \end{pmatrix} = \begin{pmatrix} 1 \\ \mathbf{x} \end{pmatrix}$$

We want to solve the independence of w_0

We would like w_0 as part of the dot product by making $x_0 = 1$

$$g(\mathbf{x}) = w_0 \times 1 + \sum_{i=1}^d w_i x_i = w_0 \times x_0 + \sum_{i=1}^d w_i x_i = \sum_{i=0}^d w_i x_i \quad (7)$$

By making

$$\mathbf{x}_{aug} = \begin{pmatrix} 1 \\ x_1 \\ \vdots \\ x_d \end{pmatrix} = \begin{pmatrix} 1 \\ \mathbf{x} \end{pmatrix}$$

where

\mathbf{x}_{aug} is called an augmented feature vector.

We want to solve the independence of w_0

We would like w_0 as part of the dot product by making $x_0 = 1$

$$g(\mathbf{x}) = w_0 \times 1 + \sum_{i=1}^d w_i x_i = w_0 \times x_0 + \sum_{i=1}^d w_i x_i = \sum_{i=0}^d w_i x_i \quad (7)$$

By making

$$\mathbf{x}_{aug} = \begin{pmatrix} 1 \\ x_1 \\ \vdots \\ x_d \end{pmatrix} = \begin{pmatrix} 1 \\ \mathbf{x} \end{pmatrix}$$

Where

\mathbf{x}_{aug} is called an augmented feature vector.

We want to solve the independence of w_0

We would like w_0 as part of the dot product by making $x_0 = 1$

$$g(\mathbf{x}) = w_0 \times 1 + \sum_{i=1}^d w_i x_i = w_0 \times x_0 + \sum_{i=1}^d w_i x_i = \sum_{i=0}^d w_i x_i \quad (7)$$

By making

$$\mathbf{x}_{aug} = \begin{pmatrix} 1 \\ x_1 \\ \vdots \\ x_d \end{pmatrix} = \begin{pmatrix} 1 \\ \mathbf{x} \end{pmatrix}$$

Where

\mathbf{x}_{aug} is called an augmented feature vector.

In a similar way

We have the augmented weight vector

$$\mathbf{w}_{aug} = \begin{pmatrix} w_0 \\ w_1 \\ \vdots \\ w_d \end{pmatrix} = \begin{pmatrix} w_0 \\ \mathbf{w} \end{pmatrix}$$



In a similar way

We have the augmented weight vector

$$\mathbf{w}_{aug} = \begin{pmatrix} w_0 \\ w_1 \\ \vdots \\ w_d \end{pmatrix} = \begin{pmatrix} w_0 \\ \mathbf{w} \end{pmatrix}$$

Remarks

- The addition of a constant component to \mathbf{x} preserves all the distance relationship between samples.
- The resulting \mathbf{x}_{aug} vectors, all lie in a d -dimensional subspace which is the \mathbf{x} -space itself.



In a similar way

We have the augmented weight vector

$$\mathbf{w}_{aug} = \begin{pmatrix} w_0 \\ w_1 \\ \vdots \\ w_d \end{pmatrix} = \begin{pmatrix} w_0 \\ \mathbf{w} \end{pmatrix}$$

Remarks

- The addition of a constant component to \mathbf{x} preserves all the distance relationship between samples.
- The resulting \mathbf{x}_{aug} vectors, all lie in a d -dimensional subspace which is the \mathbf{x} -space itself.

More Remarks

In addition

The hyperplane decision surface \hat{H} defined by

$$\mathbf{w}_{aug}^T \mathbf{x}_{aug} = 0$$

passes through the origin in \mathbf{x}_{aug} -space.

Even though

The corresponding hyperplane H can be in any position of the \mathbf{x} -space.



More Remarks

In addition

The hyperplane decision surface \hat{H} defined by

$$\mathbf{w}_{aug}^T \mathbf{x}_{aug} = 0$$

passes through the origin in \mathbf{x}_{aug} -space.

Even Though

The corresponding hyperplane H can be in any position of the \mathbf{x} -space.



More Remarks

In addition

The distance from \mathbf{y} to \hat{H} is:

$$\frac{|\mathbf{w}_{aug}^T \mathbf{x}_{aug}|}{\|\mathbf{w}_{aug}\|} = \frac{|g(\mathbf{x}_{aug})|}{\|\mathbf{w}_{aug}\|}$$



Now

Is $\|w\| \leq \|w_{aug}\|$

- Ideas?

$$\sqrt{\sum_{i=1}^d w_i^2} \leq \sqrt{\sum_{i=1}^d w_i^2 + w_0^2}$$

This mapping is quite useful

Because we only need to find a weight vector w_{aug} instead of finding the weight vector w and the threshold w_0 .



Now

Is $\|w\| \leq \|w_{aug}\|$

- Ideas?

$$\sqrt{\sum_{i=1}^d w_i^2} \leq \sqrt{\sum_{i=1}^d w_i^2 + w_0^2}$$

This mapping is quite useful

Because we only need to find a weight vector w_{aug} instead of finding the weight vector w and the threshold w_0 .



Outline

1 Linear Transformation

- Introduction
- Functions that can be defined using matrices
- Linear Functions
- Kernel and Range
- The Matrix of a Linear Transformation
- Going Back to Homogeneous Equations
- The Rank-Nullity Theorem

2 Derivative of Transformations

- Introduction
- Derivative of a Linear Transformation
- Derivative of a Quadratic Transformation

3 Linear Regression

- The Simplest Functions
- Splitting the Space
- Defining the Decision Surface
- Properties of the Hyperplane $w^T x + w_0$
- Augmenting the Vector
- **Least Squared Error Procedure**
 - The Geometry of a Two-Category Linearly-Separable Case
 - The Error Idea
 - The Final Error Equation

4 Principal Component Analysis

- Karhunen-Loeve Transform
- Projecting the Data
- Lagrange Multipliers
- The Process
- Example

5 Singular Value Decomposition

- Introduction
- Image Compression



Outline

1 Linear Transformation

- Introduction
- Functions that can be defined using matrices
- Linear Functions
- Kernel and Range
- The Matrix of a Linear Transformation
- Going Back to Homogeneous Equations
- The Rank-Nullity Theorem

2 Derivative of Transformations

- Introduction
- Derivative of a Linear Transformation
- Derivative of a Quadratic Transformation

3 Linear Regression

- The Simplest Functions
- Splitting the Space
- Defining the Decision Surface
- Properties of the Hyperplane $w^T x + w_0$
- Augmenting the Vector
- **Least Squared Error Procedure**
 - **The Geometry of a Two-Category Linearly-Separable Case**
 - The Error Idea
 - The Final Error Equation

4 Principal Component Analysis

- Karhunen-Loeve Transform
- Projecting the Data
- Lagrange Multipliers
- The Process
- Example

5 Singular Value Decomposition

- Introduction
- Image Compression



Initial Supposition

Suppose, we have

n samples $\mathbf{x}_1, \mathbf{x}_2, \dots, \mathbf{x}_n$ some labeled ω_1 and some labeled ω_2 .



Initial Supposition

Suppose, we have

n samples $\mathbf{x}_1, \mathbf{x}_2, \dots, \mathbf{x}_n$ some labeled ω_1 and some labeled ω_2 .

We want a vector weight \mathbf{w} such that

- $\mathbf{w}^T \mathbf{x}_i > 0$, if $\mathbf{x}_i \in \omega_1$.
- $\mathbf{w}^T \mathbf{x}_i < 0$, if $\mathbf{x}_i \in \omega_2$.

The name of this weight vector

It is called a separating vector or solution vector.



Initial Supposition

Suppose, we have

n samples $\mathbf{x}_1, \mathbf{x}_2, \dots, \mathbf{x}_n$ some labeled ω_1 and some labeled ω_2 .

We want a vector weight \mathbf{w} such that

- $\mathbf{w}^T \mathbf{x}_i > 0$, if $\mathbf{x}_i \in \omega_1$.
- $\mathbf{w}^T \mathbf{x}_i < 0$, if $\mathbf{x}_i \in \omega_2$.

The name of this weight vector

It is called a separating vector or solution vector.



Initial Supposition

Suppose, we have

n samples $\mathbf{x}_1, \mathbf{x}_2, \dots, \mathbf{x}_n$ some labeled ω_1 and some labeled ω_2 .

We want a vector weight \mathbf{w} such that

- $\mathbf{w}^T \mathbf{x}_i > 0$, if $\mathbf{x}_i \in \omega_1$.
- $\mathbf{w}^T \mathbf{x}_i < 0$, if $\mathbf{x}_i \in \omega_2$.

The name of this weight vector

It is called a separating vector or solution vector.



Now, assume the following

Imagine that your problem has two classes ω_1 and ω_2 in \mathbb{R}^2

① They are linearly separable!!!

② You require to label them.



Now, assume the following

Imagine that your problem has two classes ω_1 and ω_2 in \mathbb{R}^2

- 1 They are linearly separable!!!
- 2 You require to label them.

We have a problem!!!

Which is the problem?



Now, assume the following

Imagine that your problem has two classes ω_1 and ω_2 in \mathbb{R}^2

- 1 They are linearly separable!!!
- 2 You require to label them.

We have a problem!!!

Which is the problem?

We do not know the hyperplane!

Thus, what distance each point has to the hyperplane?



Now, assume the following

Imagine that your problem has two classes ω_1 and ω_2 in \mathbb{R}^2

- 1 They are linearly separable!!!
- 2 You require to label them.

We have a problem!!!

Which is the problem?

We do not know the hyperplane!!!

Thus, what distance each point has to the hyperplane?



A Simple Solution For Our Quandary

Label the Classes

- $\omega_1 \implies +1$
- $\omega_2 \implies -1$



A Simple Solution For Our Quandary

Label the Classes

- $\omega_1 \implies +1$
- $\omega_2 \implies -1$

We produce the following labels

- 1 if $\mathbf{x} \in \omega_1$ then $y_{ideal} = g_{ideal}(\mathbf{x}) = +1$.
- 2 if $\mathbf{x} \in \omega_2$ then $y_{ideal} = g_{ideal}(\mathbf{x}) = -1$.

Remark: We have a problem with this labels!!!



A Simple Solution For Our Quandary

Label the Classes

- $\omega_1 \implies +1$
- $\omega_2 \implies -1$

We produce the following labels

- 1 if $\mathbf{x} \in \omega_1$ then $y_{ideal} = g_{ideal}(\mathbf{x}) = +1$.
- 2 if $\mathbf{x} \in \omega_2$ then $y_{ideal} = g_{ideal}(\mathbf{x}) = -1$.

Remark: We have a problem with this labels!!!



A Simple Solution For Our Quandary

Label the Classes

- $\omega_1 \implies +1$
- $\omega_2 \implies -1$

We produce the following labels

- 1 if $\mathbf{x} \in \omega_1$ then $y_{ideal} = g_{ideal}(\mathbf{x}) = +1$.
- 2 if $\mathbf{x} \in \omega_2$ then $y_{ideal} = g_{ideal}(\mathbf{x}) = -1$.

Remark: We have a problem with this labels!!!



Outline

1 Linear Transformation

- Introduction
- Functions that can be defined using matrices
- Linear Functions
- Kernel and Range
- The Matrix of a Linear Transformation
- Going Back to Homogeneous Equations
- The Rank-Nullity Theorem

2 Derivative of Transformations

- Introduction
- Derivative of a Linear Transformation
- Derivative of a Quadratic Transformation

3 Linear Regression

- The Simplest Functions
- Splitting the Space
- Defining the Decision Surface
- Properties of the Hyperplane $w^T x + w_0$
- Augmenting the Vector
- Least Squared Error Procedure
- The Geometry of a Two-Category Linearly-Separable Case
- **The Error Idea**
- The Final Error Equation

4 Principal Component Analysis

- Karhunen-Loeve Transform
- Projecting the Data
- Lagrange Multipliers
- The Process
- Example

5 Singular Value Decomposition

- Introduction
- Image Compression



Now, What?

Assume true function f is given by

$$y_{noise} = g_{noise}(\mathbf{x}) = \mathbf{w}^T \mathbf{x} + w_0 + e \quad (8)$$

Where the

It has a $e \sim N(\mu, \sigma^2)$

Thus, we can do the following

$$y_{noise} = g_{noise}(\mathbf{x}) = g_{ideal}(\mathbf{x}) + e \quad (9)$$



Now, What?

Assume true function f is given by

$$y_{noise} = g_{noise}(\mathbf{x}) = \mathbf{w}^T \mathbf{x} + w_0 + e \quad (8)$$

Where the e

It has a $e \sim N(\mu, \sigma^2)$

Thus, we can do the following

$$y_{noise} = g_{noise}(\mathbf{x}) = g_{ideal}(\mathbf{x}) + e \quad (9)$$



Now, What?

Assume true function f is given by

$$y_{noise} = g_{noise}(\mathbf{x}) = \mathbf{w}^T \mathbf{x} + w_0 + e \quad (8)$$

Where the e

It has a $e \sim N(\mu, \sigma^2)$

Thus, we can do the following

$$y_{noise} = g_{noise}(\mathbf{x}) = g_{ideal}(\mathbf{x}) + e \quad (9)$$



Thus, we have

What to do?

$$e = y_{noise} - g_{ideal}(\mathbf{x}) \quad (10)$$

Graphically

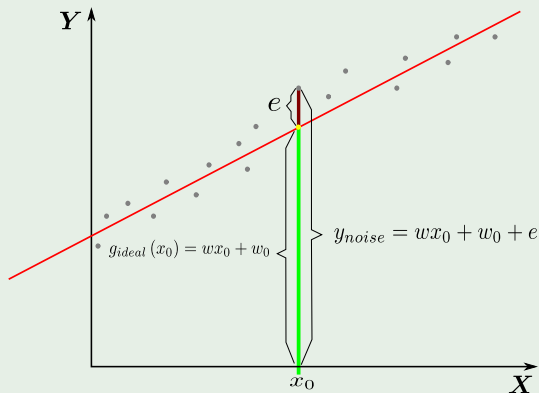


Thus, we have

What to do?

$$e = y_{\text{noise}} - g_{\text{ideal}}(\mathbf{x}) \quad (10)$$

Graphically



Then, we have

A TRICK... Quite a good one!!! Instead of using y_{noise}

$$e = y_{noise} - g_{ideal}(\mathbf{x}) \quad (11)$$

We use y_{ideal}

$$e = y_{ideal} - g_{ideal}(\mathbf{x}) \quad (12)$$

We will see

How the geometry will solve the problem with using these labels.



Then, we have

A TRICK... Quite a good one!!! Instead of using y_{noise}

$$e = y_{noise} - g_{ideal}(\mathbf{x}) \quad (11)$$

We use y_{ideal}

$$e = y_{ideal} - g_{ideal}(\mathbf{x}) \quad (12)$$

We will see

How the geometry will solve the problem with using these labels.



Then, we have

A TRICK... Quite a good one!!! Instead of using y_{noise}

$$e = y_{noise} - g_{ideal}(\mathbf{x}) \quad (11)$$

We use y_{ideal}

$$e = y_{ideal} - g_{ideal}(\mathbf{x}) \quad (12)$$

We will see

How the geometry will solve the problem with using these labels.



Outline

1 Linear Transformation

- Introduction
- Functions that can be defined using matrices
- Linear Functions
- Kernel and Range
- The Matrix of a Linear Transformation
- Going Back to Homogeneous Equations
- The Rank-Nullity Theorem

2 Derivative of Transformations

- Introduction
- Derivative of a Linear Transformation
- Derivative of a Quadratic Transformation

3 Linear Regression

- The Simplest Functions
- Splitting the Space
- Defining the Decision Surface
- Properties of the Hyperplane $w^T x + w_0$
- Augmenting the Vector
- Least Squared Error Procedure
 - The Geometry of a Two-Category Linearly-Separable Case
- The Error Idea
- **The Final Error Equation**

4 Principal Component Analysis

- Karhunen-Loeve Transform
- Projecting the Data
- Lagrange Multipliers
- The Process
- Example

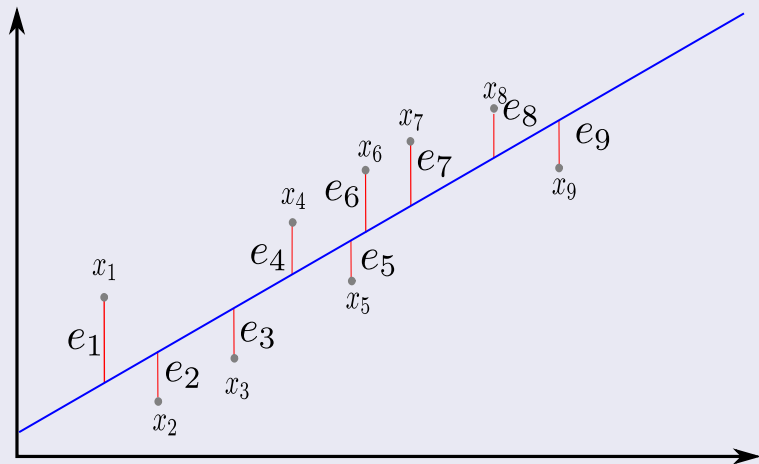
5 Singular Value Decomposition

- Introduction
- Image Compression



Here, we have multiple errors

What can we do?



Sum Over All the Errors

We can do the following

$$J(\mathbf{w}) = \sum_{i=1}^N e_i^2 = \sum_{i=1}^N (y_i - g_{ideal}(\mathbf{x}_i))^2 \quad (13)$$

Remark: This is known as the Least Squared Error cost function



Sum Over All the Errors

We can do the following

$$J(\mathbf{w}) = \sum_{i=1}^N e_i^2 = \sum_{i=1}^N (y_i - g_{ideal}(\mathbf{x}_i))^2 \quad (13)$$

Remark: This is known as the Least Squared Error cost function

Generalizing

- The dimensionality of each sample (data point) is d .

• You can extend each vector sample to be $\mathbf{x}' = (1, \mathbf{x}')$.



Sum Over All the Errors

We can do the following

$$J(\mathbf{w}) = \sum_{i=1}^N e_i^2 = \sum_{i=1}^N (y_i - g_{ideal}(\mathbf{x}_i))^2 \quad (13)$$

Remark: This is known as the Least Squared Error cost function

Generalizing

- The dimensionality of each sample (data point) is d .
- You can extend each vector sample to be $\mathbf{x}^T = (\mathbf{1}, \mathbf{x}')$.



We can use a trick

The following function

$$g_{ideal}(\mathbf{x}) = \begin{pmatrix} 1 & x_1 & x_2 & \dots & x_d \end{pmatrix} \begin{pmatrix} w_0 \\ w_1 \\ w_2 \\ w_3 \\ \vdots \\ w_d \end{pmatrix} = \mathbf{x}^T \mathbf{w}$$

We can rewrite the error equation as

$$J(\mathbf{w}) = \sum_{i=1}^N (y_i - g_{ideal}(\mathbf{x}_i))^2 = \sum_{i=1}^N (y_i - \mathbf{x}_i^T \mathbf{w})^2 \quad (14)$$



We can use a trick

The following function

$$g_{ideal}(\mathbf{x}) = \begin{pmatrix} 1 & x_1 & x_2 & \dots & x_d \end{pmatrix} \begin{pmatrix} w_0 \\ w_1 \\ w_2 \\ w_3 \\ \vdots \\ w_d \end{pmatrix} = \mathbf{x}^T \mathbf{w}$$

We can rewrite the error equation as

$$J(\mathbf{w}) = \sum_{i=1}^N (y_i - g_{ideal}(\mathbf{x}_i))^2 = \sum_{i=1}^N (y_i - \mathbf{x}_i^T \mathbf{w})^2 \quad (14)$$



Furthermore

Then stacking all the possible estimations into the product Data Matrix and weight vector

$$\mathbf{X}\mathbf{w} = \begin{pmatrix} 1 & (\mathbf{x}_1)_1 & \cdots & (\mathbf{x}_1)_j & \cdots & (\mathbf{x}_1)_d \\ \vdots & & & \vdots & & \vdots \\ 1 & (\mathbf{x}_i)_1 & & (\mathbf{x}_i)_j & & (\mathbf{x}_i)_d \\ \vdots & & & \vdots & & \vdots \\ 1 & (\mathbf{x}_N)_1 & \cdots & (\mathbf{x}_N)_j & \cdots & (\mathbf{x}_N)_d \end{pmatrix} \begin{pmatrix} w_1 \\ w_2 \\ w_3 \\ \vdots \\ w_{d+1} \end{pmatrix}$$



Note about other representations

We could have $\mathbf{x}^T = (x_1, x_2, \dots, x_d, 1)$ thus

$$\mathbf{X} = \begin{pmatrix} (\mathbf{x}_1)_1 & \cdots & (\mathbf{x}_1)_j & \cdots & (\mathbf{x}_1)_d & 1 \\ & & \vdots & & \vdots & \vdots \\ (\mathbf{x}_i)_1 & & (\mathbf{x}_i)_j & & (\mathbf{x}_i)_d & 1 \\ & & \vdots & & \vdots & \vdots \\ (\mathbf{x}_N)_1 & \cdots & (\mathbf{x}_N)_j & \cdots & (\mathbf{x}_N)_d & 1 \end{pmatrix} \quad (15)$$



Then, we have the following trick with \mathbf{X}

With the Data Matrix

$$\mathbf{X}w = \begin{pmatrix} \mathbf{x}_1^T w \\ \mathbf{x}_2^T w \\ \mathbf{x}_3^T w \\ \vdots \\ \mathbf{x}_N^T w \end{pmatrix} \quad (16)$$



Therefore

We have that

$$\begin{pmatrix} y_1 \\ y_2 \\ y_3 \\ \vdots \\ y_4 \end{pmatrix} - \begin{pmatrix} \mathbf{x}_1^T \mathbf{w} \\ \mathbf{x}_2^T \mathbf{w} \\ \mathbf{x}_3^T \mathbf{w} \\ \vdots \\ \mathbf{x}_N^T \mathbf{w} \end{pmatrix} = \begin{pmatrix} y_1 - \mathbf{x}_1^T \mathbf{w} \\ y_2 - \mathbf{x}_2^T \mathbf{w} \\ y_3 - \mathbf{x}_3^T \mathbf{w} \\ \vdots \\ y_4 - \mathbf{x}_N^T \mathbf{w} \end{pmatrix}$$

Then, we have the following equality:

$$\begin{pmatrix} y_1 - \mathbf{x}_1^T \mathbf{w} & y_2 - \mathbf{x}_2^T \mathbf{w} & y_3 - \mathbf{x}_3^T \mathbf{w} & \dots & y_4 - \mathbf{x}_N^T \mathbf{w} \end{pmatrix} \begin{pmatrix} y_1 - \mathbf{x}_1^T \mathbf{w} \\ y_2 - \mathbf{x}_2^T \mathbf{w} \\ y_3 - \mathbf{x}_3^T \mathbf{w} \\ \vdots \\ y_4 - \mathbf{x}_N^T \mathbf{w} \end{pmatrix} = \sum_{i=1}^N (y_i - \mathbf{x}_i^T \mathbf{w})^2$$

Therefore

We have that

$$\begin{pmatrix} y_1 \\ y_2 \\ y_3 \\ \vdots \\ y_4 \end{pmatrix} - \begin{pmatrix} \mathbf{x}_1^T \mathbf{w} \\ \mathbf{x}_2^T \mathbf{w} \\ \mathbf{x}_3^T \mathbf{w} \\ \vdots \\ \mathbf{x}_N^T \mathbf{w} \end{pmatrix} = \begin{pmatrix} y_1 - \mathbf{x}_1^T \mathbf{w} \\ y_2 - \mathbf{x}_2^T \mathbf{w} \\ y_3 - \mathbf{x}_3^T \mathbf{w} \\ \vdots \\ y_4 - \mathbf{x}_N^T \mathbf{w} \end{pmatrix}$$

Then, we have the following equality

$$\begin{pmatrix} y_1 - \mathbf{x}_1^T \mathbf{w} & y_2 - \mathbf{x}_2^T \mathbf{w} & y_3 - \mathbf{x}_3^T \mathbf{w} & \cdots & y_4 - \mathbf{x}_N^T \mathbf{w} \end{pmatrix} \begin{pmatrix} y_1 - \mathbf{x}_1^T \mathbf{w} \\ y_2 - \mathbf{x}_2^T \mathbf{w} \\ y_3 - \mathbf{x}_3^T \mathbf{w} \\ \vdots \\ y_4 - \mathbf{x}_N^T \mathbf{w} \end{pmatrix} = \sum_{i=1}^N (y_i - \mathbf{x}_i^T \mathbf{w})^2$$

Then, we have

The following equality

$$\sum_{i=1}^N \left(y_i - \mathbf{x}_i^T \mathbf{w} \right)^2 = (\mathbf{y} - \mathbf{X}\mathbf{w})^T (\mathbf{y} - \mathbf{X}\mathbf{w}) = \|\mathbf{y} - \mathbf{X}\mathbf{w}\|_2^2 \quad (17)$$



We can expand our quadratic formula!!!

Thus

$$(\mathbf{y} - \mathbf{X}\mathbf{w})^T (\mathbf{y} - \mathbf{X}\mathbf{w}) = \mathbf{y}^T \mathbf{y} - \mathbf{y}^T \mathbf{X}\mathbf{w} - \mathbf{w}^T \mathbf{X}^T \mathbf{y} + \mathbf{w}^T \mathbf{X}^T \mathbf{X}\mathbf{w} \quad (18)$$



We can expand our quadratic formula!!!

Thus

$$(\mathbf{y} - \mathbf{X}\mathbf{w})^T (\mathbf{y} - \mathbf{X}\mathbf{w}) = \mathbf{y}^T \mathbf{y} - \mathbf{y}^T \mathbf{X}\mathbf{w} - \mathbf{w}^T \mathbf{X}^T \mathbf{y} + \mathbf{w}^T \mathbf{X}^T \mathbf{X}\mathbf{w} \quad (18)$$

Now

- Derive with respect to \mathbf{w}

• Assume that $\mathbf{X}^T \mathbf{X}$ is invertible



We can expand our quadratic formula!!!

Thus

$$(\mathbf{y} - \mathbf{X}\mathbf{w})^T (\mathbf{y} - \mathbf{X}\mathbf{w}) = \mathbf{y}^T \mathbf{y} - \mathbf{y}^T \mathbf{X}\mathbf{w} - \mathbf{w}^T \mathbf{X}^T \mathbf{y} + \mathbf{w}^T \mathbf{X}^T \mathbf{X}\mathbf{w} \quad (18)$$

Now

- Derive with respect to \mathbf{w}
- Assume that $\mathbf{X}^T \mathbf{X}$ is invertible



Therefore

We have the following equivalences

$$\frac{d\mathbf{w}^T A \mathbf{w}}{d\mathbf{w}} = \mathbf{w}^T (A + A^T), \quad \frac{d\mathbf{w}^T A}{d\mathbf{w}} = A^T \quad (19)$$

Now given that the transpose of a number is the number itself

$$y^T X w = [y^T X w]^T = w^T X^T y$$



Therefore

We have the following equivalences

$$\frac{d\mathbf{w}^T A \mathbf{w}}{d\mathbf{w}} = \mathbf{w}^T (A + A^T), \quad \frac{d\mathbf{w}^T A}{d\mathbf{w}} = A^T \quad (19)$$

Now given that the transpose of a number is the number itself

$$\mathbf{y}^T \mathbf{X} \mathbf{w} = [\mathbf{y}^T \mathbf{X} \mathbf{w}]^T = \mathbf{w}^T \mathbf{X}^T \mathbf{y}$$



Then, when we derive by w

We have then

$$\frac{d \left(\mathbf{y}^T \mathbf{y} - 2\mathbf{w}^T \mathbf{X}^T \mathbf{y} + \mathbf{w}^T \mathbf{X}^T \mathbf{X} \mathbf{w} \right)}{d\mathbf{w}} = -2\mathbf{y}^T \mathbf{X} + \mathbf{w}^T \left(\mathbf{X}^T \mathbf{X} + \left(\mathbf{X}^T \mathbf{X} \right) \right)$$
$$= -2\mathbf{y}^T \mathbf{X} + 2\mathbf{w}^T \left(\mathbf{X}^T \mathbf{X} \right)$$

Then, when we derive by w

We have then

$$\begin{aligned}\frac{d\left(\mathbf{y}^T \mathbf{y} - 2\mathbf{w}^T \mathbf{X}^T \mathbf{y} + \mathbf{w}^T \mathbf{X}^T \mathbf{X} \mathbf{w}\right)}{d\mathbf{w}} &= -2\mathbf{y}^T \mathbf{X} + \mathbf{w}^T \left(\mathbf{X}^T \mathbf{X} + \left(\mathbf{X}^T \mathbf{X}\right)\right) \\ &= -2\mathbf{y}^T \mathbf{X} + 2\mathbf{w}^T \left(\mathbf{X}^T \mathbf{X}\right)\end{aligned}$$

Making this equal to the zero row vector

$$-2\mathbf{y}^T \mathbf{X} + 2\mathbf{w}^T \left(\mathbf{X}^T \mathbf{X}\right) = 0$$

Then, when we derive by w

We have then

$$\begin{aligned}\frac{d\left(\mathbf{y}^T \mathbf{y} - 2\mathbf{w}^T \mathbf{X}^T \mathbf{y} + \mathbf{w}^T \mathbf{X}^T \mathbf{X} \mathbf{w}\right)}{d\mathbf{w}} &= -2\mathbf{y}^T \mathbf{X} + \mathbf{w}^T \left(\mathbf{X}^T \mathbf{X} + \left(\mathbf{X}^T \mathbf{X}\right)\right) \\ &= -2\mathbf{y}^T \mathbf{X} + 2\mathbf{w}^T \left(\mathbf{X}^T \mathbf{X}\right)\end{aligned}$$

Making this equal to the zero row vector

$$-2\mathbf{y}^T \mathbf{X} + 2\mathbf{w}^T \left(\mathbf{X}^T \mathbf{X}\right) = 0$$

We apply the transpose

$$\begin{aligned}\left[-2\mathbf{y}^T \mathbf{X} + 2\mathbf{w}^T \left(\mathbf{X}^T \mathbf{X}\right)\right]^T &= [0]^T \\ -2\mathbf{X}^T \mathbf{y} + 2\left(\mathbf{X}^T \mathbf{X}\right) \mathbf{w} &= 0 \text{ (column vector)}\end{aligned}$$

Then, when we derive by w

We have then

$$\begin{aligned}\frac{d\left(\mathbf{y}^T \mathbf{y} - 2\mathbf{w}^T \mathbf{X}^T \mathbf{y} + \mathbf{w}^T \mathbf{X}^T \mathbf{X} \mathbf{w}\right)}{d\mathbf{w}} &= -2\mathbf{y}^T \mathbf{X} + \mathbf{w}^T \left(\mathbf{X}^T \mathbf{X} + \left(\mathbf{X}^T \mathbf{X}\right)\right) \\ &= -2\mathbf{y}^T \mathbf{X} + 2\mathbf{w}^T \left(\mathbf{X}^T \mathbf{X}\right)\end{aligned}$$

Making this equal to the zero row vector

$$-2\mathbf{y}^T \mathbf{X} + 2\mathbf{w}^T \left(\mathbf{X}^T \mathbf{X}\right) = 0$$

We apply the transpose

$$\left[-2\mathbf{y}^T \mathbf{X} + 2\mathbf{w}^T \left(\mathbf{X}^T \mathbf{X}\right)\right]^T = [0]^T$$

$$-2\mathbf{X}^T \mathbf{y} + 2\left(\mathbf{X}^T \mathbf{X}\right) \mathbf{w} = 0 \text{ (column vector)}$$

Then, when we derive by w

We have then

$$\begin{aligned}\frac{d\left(\mathbf{y}^T \mathbf{y} - 2\mathbf{w}^T \mathbf{X}^T \mathbf{y} + \mathbf{w}^T \mathbf{X}^T \mathbf{X} \mathbf{w}\right)}{d\mathbf{w}} &= -2\mathbf{y}^T \mathbf{X} + \mathbf{w}^T \left(\mathbf{X}^T \mathbf{X} + \left(\mathbf{X}^T \mathbf{X}\right)\right) \\ &= -2\mathbf{y}^T \mathbf{X} + 2\mathbf{w}^T \left(\mathbf{X}^T \mathbf{X}\right)\end{aligned}$$

Making this equal to the zero row vector

$$-2\mathbf{y}^T \mathbf{X} + 2\mathbf{w}^T \left(\mathbf{X}^T \mathbf{X}\right) = 0$$

We apply the transpose

$$\begin{aligned}\left[-2\mathbf{y}^T \mathbf{X} + 2\mathbf{w}^T \left(\mathbf{X}^T \mathbf{X}\right)\right]^T &= [0]^T \\ -2\mathbf{X}^T \mathbf{y} + 2\left(\mathbf{X}^T \mathbf{X}\right) \mathbf{w} &= 0 \text{ (column vector)}\end{aligned}$$

Solving for w

We have then

$$w = (\mathbf{X}^T \mathbf{X})^{-1} \mathbf{X}^T \mathbf{y} \quad (20)$$

Note: $\mathbf{X}^T \mathbf{X}$ is always positive semi-definite. If it is also invertible, it is positive definite.

Hint: How we get the discriminant function?

Any Ideas?



Solving for w

We have then

$$w = \left(X^T X \right)^{-1} X^T y \quad (20)$$

Note: $X^T X$ is always positive semi-definite. If it is also invertible, it is positive definite.

Thus, How we get the discriminant function?

Any Ideas?



The Final Discriminant Function

Very Simple!!!

$$g(\mathbf{x}) = \mathbf{x}^T \mathbf{w} = \mathbf{x}^T (\mathbf{X}^T \mathbf{X})^{-1} \mathbf{X}^T \mathbf{y} \quad (21)$$



Outline

1 Linear Transformation

- Introduction
- Functions that can be defined using matrices
- Linear Functions
- Kernel and Range
- The Matrix of a Linear Transformation
- Going Back to Homogeneous Equations
- The Rank-Nullity Theorem

2 Derivative of Transformations

- Introduction
- Derivative of a Linear Transformation
- Derivative of a Quadratic Transformation

3 Linear Regression

- The Simplest Functions
- Splitting the Space
- Defining the Decision Surface
- Properties of the Hyperplane $w^T x + w_0$
- Augmenting the Vector
- Least Squared Error Procedure
 - The Geometry of a Two-Category Linearly-Separable Case
- The Error Idea
- The Final Error Equation

4 Principal Component Analysis

- **Karhunen-Loeve Transform**
- Projecting the Data
- Lagrange Multipliers
- The Process
- Example

5 Singular Value Decomposition

- Introduction
- Image Compression

Also Known as Karhunen-Loeve Transform

Setup

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

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 $\{\mathbf{x}_n\}$ with $n = 1, 2, \dots, N$ and $\mathbf{x}_n \in \mathbb{R}^d$.

Goal

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



Dimensional Variance

Remember the Variance Sample in \mathbb{R}

$$VAR(X) = \frac{\sum_{i=1}^N (x_i - \bar{x})^2}{N - 1} \quad (22)$$

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

$$COV(x, y) = \frac{\sum_{i=1}^N (x_i - \bar{x})(y_i - \bar{y})}{N - 1} \quad (23)$$



Dimensional Variance

Remember the Variance Sample in \mathbb{R}

$$VAR(X) = \frac{\sum_{i=1}^N (x_i - \bar{x})^2}{N - 1} \quad (22)$$

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

$$COV(x, y) = \frac{\sum_{i=1}^N (x_i - \bar{x})(y_i - \bar{y})}{N - 1} \quad (23)$$



Now, Define

Given the data

$$\mathbf{x}_1, \mathbf{x}_2, \dots, \mathbf{x}_N \quad (24)$$

where \mathbf{x}_i is a column vector

Construct the sample mean

$$\bar{\mathbf{x}} = \frac{1}{N} \sum_{i=1}^N \mathbf{x}_i \quad (25)$$

Center data

$$\mathbf{x}_1 - \bar{\mathbf{x}}, \mathbf{x}_2 - \bar{\mathbf{x}}, \dots, \mathbf{x}_N - \bar{\mathbf{x}} \quad (26)$$



Now, Define

Given the data

$$\mathbf{x}_1, \mathbf{x}_2, \dots, \mathbf{x}_N \quad (24)$$

where \mathbf{x}_i is a column vector

Construct the sample mean

$$\bar{\mathbf{x}} = \frac{1}{N} \sum_{i=1}^N \mathbf{x}_i \quad (25)$$

Center data

$$\mathbf{x}_1 - \bar{\mathbf{x}}, \mathbf{x}_2 - \bar{\mathbf{x}}, \dots, \mathbf{x}_N - \bar{\mathbf{x}} \quad (26)$$



Now, Define

Given the data

$$\mathbf{x}_1, \mathbf{x}_2, \dots, \mathbf{x}_N \quad (24)$$

where \mathbf{x}_i is a column vector

Construct the sample mean

$$\bar{\mathbf{x}} = \frac{1}{N} \sum_{i=1}^N \mathbf{x}_i \quad (25)$$

Center data

$$\mathbf{x}_1 - \bar{\mathbf{x}}, \mathbf{x}_2 - \bar{\mathbf{x}}, \dots, \mathbf{x}_N - \bar{\mathbf{x}} \quad (26)$$



Build the Sample Mean

The Covariance Matrix

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

Properties

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



Build the Sample Mean

The Covariance Matrix

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

Properties

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



Outline

1 Linear Transformation

- Introduction
- Functions that can be defined using matrices
- Linear Functions
- Kernel and Range
- The Matrix of a Linear Transformation
- Going Back to Homogeneous Equations
- The Rank-Nullity Theorem

2 Derivative of Transformations

- Introduction
- Derivative of a Linear Transformation
- Derivative of a Quadratic Transformation

3 Linear Regression

- The Simplest Functions
- Splitting the Space
- Defining the Decision Surface
- Properties of the Hyperplane $w^T x + w_0$
- Augmenting the Vector
- Least Squared Error Procedure
 - The Geometry of a Two-Category Linearly-Separable Case
- The Error Idea
- The Final Error Equation

4 Principal Component Analysis

- Karhunen-Loeve Transform
- **Projecting the Data**
- Lagrange Multipliers
- The Process
- Example

5 Singular Value Decomposition

- Introduction
- Image Compression



Using S to Project Data

For this we use a u_1

- with $u_1^T u_1 = 1$, 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?



Cinvestav

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 Linear Transformation

- Introduction
- Functions that can be defined using matrices
- Linear Functions
- Kernel and Range
- The Matrix of a Linear Transformation
- Going Back to Homogeneous Equations
- The Rank-Nullity Theorem

2 Derivative of Transformations

- Introduction
- Derivative of a Linear Transformation
- Derivative of a Quadratic Transformation

3 Linear Regression

- The Simplest Functions
- Splitting the Space
- Defining the Decision Surface
- Properties of the Hyperplane $w^T x + w_0$
- Augmenting the Vector
- Least Squared Error Procedure
 - The Geometry of a Two-Category Linearly-Separable Case
- The Error Idea
- The Final Error Equation

4 Principal Component Analysis

- Karhunen-Loeve Transform
- Projecting the Data
- **Lagrange Multipliers**
- The Process
- Example

5 Singular Value Decomposition

- Introduction
- Image Compression



Thus we have

Variance of the projected data

$$\frac{1}{N-1} \sum_{i=1}^N [\mathbf{u}_1 \mathbf{x}_i - \mathbf{u}_1 \bar{\mathbf{x}}] = \mathbf{u}_1^T S \mathbf{u}_1 \quad (28)$$

Use Lagrange Multipliers to Maximize

$$\mathbf{u}_1^T S \mathbf{u}_1 + \lambda_1 (1 - \mathbf{u}_1^T \mathbf{u}_1) \quad (29)$$



Thus we have

Variance of the projected data

$$\frac{1}{N-1} \sum_{i=1}^N [\mathbf{u}_1 \mathbf{x}_i - \mathbf{u}_1 \bar{\mathbf{x}}] = \mathbf{u}_1^T S \mathbf{u}_1 \quad (28)$$

Use Lagrange Multipliers to Maximize

$$\mathbf{u}_1^T S \mathbf{u}_1 + \lambda_1 (1 - \mathbf{u}_1^T \mathbf{u}_1) \quad (29)$$



Derive by \mathbf{u}_1

We get

$$S\mathbf{u}_1 = \lambda_1\mathbf{u}_1 \quad (30)$$

Then

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

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

$$\mathbf{u}_1^T S \mathbf{u}_1 = \lambda_1 \quad (31)$$



Derive by \mathbf{u}_1

We get

$$S\mathbf{u}_1 = \lambda_1\mathbf{u}_1 \quad (30)$$

Then

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

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

$$\mathbf{u}_1^T S \mathbf{u}_1 = \lambda_1 \quad (31)$$



Derive by \mathbf{u}_1

We get

$$S\mathbf{u}_1 = \lambda_1\mathbf{u}_1 \quad (30)$$

Then

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

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

$$\mathbf{u}_1^T S\mathbf{u}_1 = \lambda_1 \quad (31)$$



What about the second eigenvector \mathbf{u}_2

We have the following optimization problem

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

Lagrange multiplier

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



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

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



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

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



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

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



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

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

implying

- \mathbf{u}_2 is the eigenvector of S with second largest eigenvalue λ_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 λ_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 λ_2 .



Thus

Variance will be the maximum when

$$\mathbf{u}_1^T S \mathbf{u}_1 = \lambda_1 \quad (32)$$

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

Eigen 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

$$\mathbf{u}_1^T \mathbf{S} \mathbf{u}_1 = \lambda_1 \quad (32)$$

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

- 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

$$\mathbf{u}_1^T S \mathbf{u}_1 = \lambda_1 \quad (32)$$

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

- 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 Linear Transformation

- Introduction
- Functions that can be defined using matrices
- Linear Functions
- Kernel and Range
- The Matrix of a Linear Transformation
- Going Back to Homogeneous Equations
- The Rank-Nullity Theorem

2 Derivative of Transformations

- Introduction
- Derivative of a Linear Transformation
- Derivative of a Quadratic Transformation

3 Linear Regression

- The Simplest Functions
- Splitting the Space
- Defining the Decision Surface
- Properties of the Hyperplane $w^T x + w_0$
- Augmenting the Vector
- Least Squared Error Procedure
 - The Geometry of a Two-Category Linearly-Separable Case
- The Error Idea
- The Final Error Equation

4 Principal Component Analysis

- Karhunen-Loeve Transform
- Projecting the Data
- Lagrange Multipliers
- **The Process**
- Example

5 Singular Value Decomposition

- Introduction
- Image Compression



We have the following steps

Determine covariance matrix

$$S = \frac{1}{N-1} \sum_{i=1}^N (\mathbf{x}_i - \bar{\mathbf{x}}) (\mathbf{x}_i - \bar{\mathbf{x}})^T \quad (33)$$

Generate the decomposition

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

With

- Eigenvalues in Σ and eigenvectors in the columns of U .



We have the following steps

Determine covariance matrix

$$S = \frac{1}{N-1} \sum_{i=1}^N (\mathbf{x}_i - \bar{\mathbf{x}}) (\mathbf{x}_i - \bar{\mathbf{x}})^T \quad (33)$$

Generate the decomposition

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

with

- Eigenvalues in Σ and eigenvectors in the columns of U .



We have the following steps

Determine covariance matrix

$$S = \frac{1}{N-1} \sum_{i=1}^N (\mathbf{x}_i - \bar{\mathbf{x}}) (\mathbf{x}_i - \bar{\mathbf{x}})^T \quad (33)$$

Generate the decomposition

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

With

- Eigenvalues in Σ and eigenvectors in the columns of U .



Then

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

$$z_i = U_K^T \mathbf{x}_i$$

- With U_k is a matrix with k columns



Outline

1 Linear Transformation

- Introduction
- Functions that can be defined using matrices
- Linear Functions
- Kernel and Range
- The Matrix of a Linear Transformation
- Going Back to Homogeneous Equations
- The Rank-Nullity Theorem

2 Derivative of Transformations

- Introduction
- Derivative of a Linear Transformation
- Derivative of a Quadratic Transformation

3 Linear Regression

- The Simplest Functions
- Splitting the Space
- Defining the Decision Surface
- Properties of the Hyperplane $w^T x + w_0$
- Augmenting the Vector
- Least Squared Error Procedure
 - The Geometry of a Two-Category Linearly-Separable Case
- The Error Idea
- The Final Error Equation

4 Principal Component Analysis

- Karhunen-Loeve Transform
- Projecting the Data
- Lagrange Multipliers
- The Process
- **Example**

5 Singular Value Decomposition

- Introduction
- Image Compression



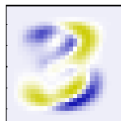
Example

From Bishop

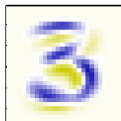
Mean



$\lambda_1 = 3.4 \cdot 10^5$



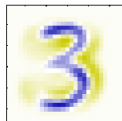
$\lambda_2 = 2.8 \cdot 10^5$



$\lambda_3 = 2.4 \cdot 10^5$



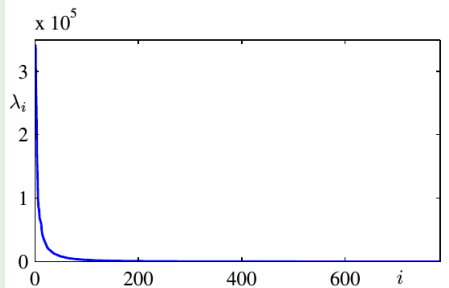
$\lambda_4 = 1.6 \cdot 10^5$



Cinvestav

Example

From Bishop



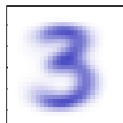
Example

From Bishop

Original



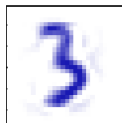
$M = 1$



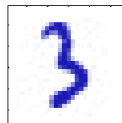
$M = 10$



$M = 50$



$M = 250$



Outline

1 Linear Transformation

- Introduction
- Functions that can be defined using matrices
- Linear Functions
- Kernel and Range
- The Matrix of a Linear Transformation
- Going Back to Homogeneous Equations
- The Rank-Nullity Theorem

2 Derivative of Transformations

- Introduction
- Derivative of a Linear Transformation
- Derivative of a Quadratic Transformation

3 Linear Regression

- The Simplest Functions
- Splitting the Space
- Defining the Decision Surface
- Properties of the Hyperplane $w^T x + w_0$
- Augmenting the Vector
- Least Squared Error Procedure
 - The Geometry of a Two-Category Linearly-Separable Case
- The Error Idea
- The Final Error Equation

4 Principal Component Analysis

- Karhunen-Loeve Transform
- Projecting the Data
- Lagrange Multipliers
- The Process
- Example

5 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 projector 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

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

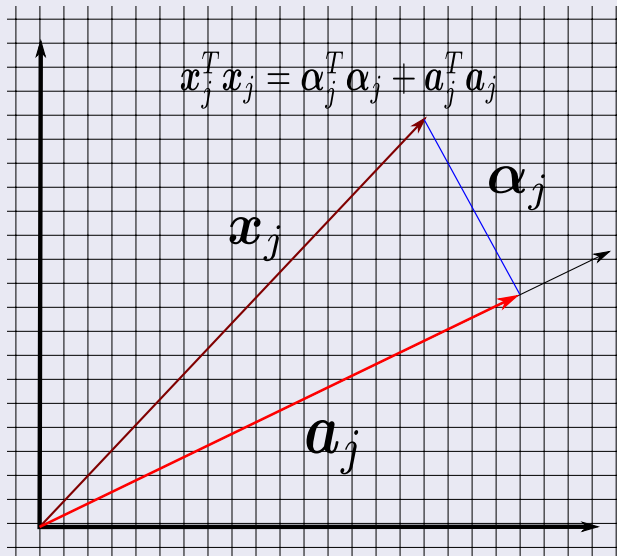
We want to know a little bit of the relations between them

- 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

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

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

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

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

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

Actually this is know 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 know 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 know 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 know 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



Cinvestav

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



Cinvestav

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 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

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

Therefore the length of the projected vector is

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



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

Definition of the projection under unitary vector

$$p = \frac{\mathbf{v}^T \mathbf{a}_i}{\mathbf{v}^T \mathbf{v}} \mathbf{v} = \left[\mathbf{v}^T \mathbf{a}_i \right] \mathbf{v}$$

Therefore the length of the projected vector is

$$\left\| \left[\mathbf{v}^T \mathbf{a}_i \right] \mathbf{v} \right\| = \left| \mathbf{v}^T \mathbf{a}_i \right|$$



Then

Thus with a little bit of notation

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

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)

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

Then, we can define the following singular value

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



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....

Conclusion

- 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 \mathbf{v}_2

$$\mathbf{v}_2 = \arg \max_{\mathbf{v} \perp \mathbf{v}_1, \|\mathbf{v}\|=1} \|A\mathbf{v}\|$$

Then you go through this process

- Stop when we have found all the following vectors:

$$\mathbf{v}_1, \mathbf{v}_2, \dots, \mathbf{v}_r$$

As singular vectors and

$$\arg \max_{\substack{\mathbf{v} \perp \mathbf{v}_1, \mathbf{v}_2, \dots, \mathbf{v}_r \\ \|\mathbf{v}\|=1}} \|A\mathbf{v}\|$$

Then, in a Greedy Fashion

The second singular vector \mathbf{v}_2

$$\mathbf{v}_2 = \arg \max_{\mathbf{v} \perp \mathbf{v}_1, \|\mathbf{v}\|=1} \|A\mathbf{v}\|$$

Then you go through this process

- Stop when we have found all the following vectors:

$$\mathbf{v}_1, \mathbf{v}_2, \dots, \mathbf{v}_r$$

As singular vectors and

$$\arg \max_{\substack{\mathbf{v} \perp \mathbf{v}_1, \mathbf{v}_2, \dots, \mathbf{v}_r \\ \|\mathbf{v}\| = 1}} \|A\mathbf{v}\|$$

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 $w_1, w_2 \in W$

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

Now, choose a basis v_1, v_2 so that v_2 is perpendicular to v_1 .

- 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 \mathbf{v}_2

- \mathbf{V}_k is at least as good as \mathbf{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

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

- 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.



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 Σ 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 Σ 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 Linear Transformation

- Introduction
- Functions that can be defined using matrices
- Linear Functions
- Kernel and Range
- The Matrix of a Linear Transformation
- Going Back to Homogeneous Equations
- The Rank-Nullity Theorem

2 Derivative of Transformations

- Introduction
- Derivative of a Linear Transformation
- Derivative of a Quadratic Transformation

3 Linear Regression

- The Simplest Functions
- Splitting the Space
- Defining the Decision Surface
- Properties of the Hyperplane $w^T x + w_0$
- Augmenting the Vector
- Least Squared Error Procedure
 - The Geometry of a Two-Category Linearly-Separable Case
- The Error Idea
- The Final Error Equation

4 Principal Component Analysis

- Karhunen-Loeve Transform
- Projecting the Data
- Lagrange Multipliers
- The Process
- Example

5 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

$$A = A_1 + A_2 + \dots + A_R \quad (34)$$

Why?

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



Singular Value Decomposition as Sums

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

$$A = A_1 + A_2 + \dots + A_R \quad (34)$$

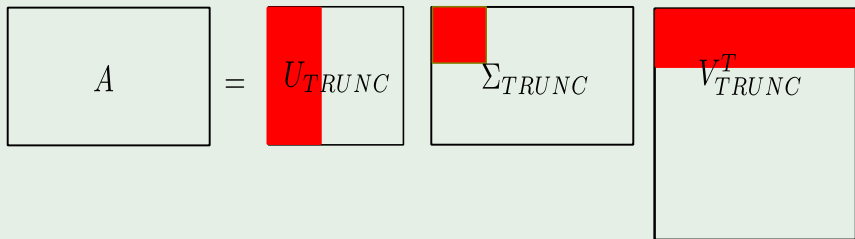
Why?

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



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 512×512

- 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 512×512

- 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 512×512

- 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 $512 \times 80 + 80 + 80 \times 512 = 82,000$