

Introduction to Machine Learning

Introduction to Linear Classifiers

Andres Mendez-Vazquez

May 24, 2020

Outline

1

Introduction

- Introduction
- Regression as approximation
- The Simplest Functions
- Splitting the Space
- Defining the Decision Surface
- Properties of the Hyperplane $w^T x + w_0$
- Augmenting the Vector

2

Developing a Solution

- Least Squared Error Procedure
 - The Geometry of a Two-Category Linearly-Separable Case
- The Error Idea
- The Final Error Equation
- Basic Solution
- Multidimensional Solution
- Remember in matrices of 3×3
- What Lives Where?
- Geometric Interpretation
- Solving the Labeling Issue
- Multi-Class Solution
- Issues with Least Squares!!!
 - Singularity Notes
 - Problem with Outliers
 - Problem with High Number of Dimensions
- What can be done?
 - Using Statistics to find Important Features
 - What about Numerical Stability?
 - Ridge Regression
- Observation About Eigenvalues

3

Exercises

- Some Stuff for the Lab



Outline

1

Introduction

● Introduction

- Regression as approximation
- The Simplest Functions
- Splitting the Space
- Defining the Decision Surface
- Properties of the Hyperplane $w^T x + w_0$
- Augmenting the Vector

2

Developing a Solution

- Least Squared Error Procedure
 - The Geometry of a Two-Category Linearly-Separable Case
- The Error Idea
- The Final Error Equation
- Basic Solution
- Multidimensional Solution
- Remember in matrices of 3×3
- What Lives Where?
- Geometric Interpretation
- Solving the Labeling Issue
- Multi-Class Solution
- Issues with Least Squares!!!
 - Singularity Notes
 - Problem with Outliers
 - Problem with High Number of Dimensions
- What can be done?
 - Using Statistics to find Important Features
 - What about Numerical Stability?
 - Ridge Regression
- Observation About Eigenvalues

3

Exercises

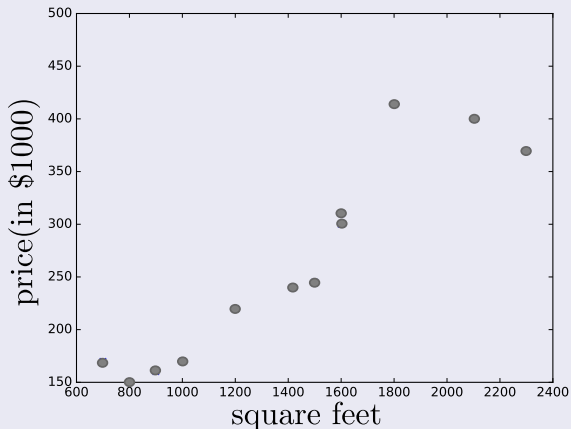
- Some Stuff for the Lab



Many Times, we have things as regression

We have this kind of data sets (House Prices per Square Feet)

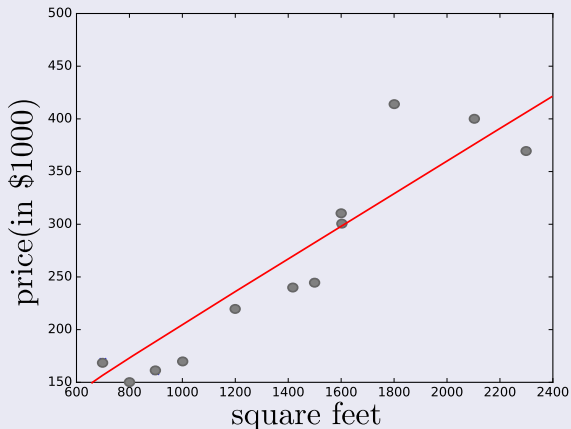
$$\begin{pmatrix} \text{Squared Feet} \\ \text{Price} \end{pmatrix} \rightarrow \begin{pmatrix} 2104 \\ 400 \end{pmatrix} \begin{pmatrix} 1800 \\ 460 \end{pmatrix} \begin{pmatrix} 1600 \\ 300 \end{pmatrix} \begin{pmatrix} 2300 \\ 370 \end{pmatrix} \dots$$



Thus

We can adjust a line/hyperplane to be able to forecast prices

$$\begin{pmatrix} \text{Squared Feet} \\ \text{Price} \end{pmatrix} \rightarrow \begin{pmatrix} 2104 \\ 400 \end{pmatrix} \begin{pmatrix} 1800 \\ 460 \end{pmatrix} \begin{pmatrix} 1600 \\ 300 \end{pmatrix} \begin{pmatrix} 2300 \\ 370 \end{pmatrix} \dots$$



Thus, Our Objective

To find such hyperplane

To do forecasting on the prices of a house given its surface!!!

Here, where Learning Machine Learning style comes around

Basically, the process defined in Machine Learning!!!



Thus, Our Objective

To find such hyperplane

To do forecasting on the prices of a house given its surface!!!

Here, where “Learning” Machine Learning style comes around

Basically, the process defined in Machine Learning!!!



Outline

1

Introduction

- Introduction
- **Regression as approximation**
- The Simplest Functions
- Splitting the Space
- Defining the Decision Surface
- Properties of the Hyperplane $w^T x + w_0$
- Augmenting the Vector

2

Developing a Solution

- Least Squared Error Procedure
 - The Geometry of a Two-Category Linearly-Separable Case
- The Error Idea
- The Final Error Equation
- Basic Solution
- Multidimensional Solution
- Remember in matrices of 3×3
- What Lives Where?
- Geometric Interpretation
- Solving the Labeling Issue
- Multi-Class Solution
- Issues with Least Squares!!!
 - Singularity Notes
 - Problem with Outliers
 - Problem with High Number of Dimensions
- What can be done?
 - Using Statistics to find Important Features
 - What about Numerical Stability?
 - Ridge Regression
- Observation About Eigenvalues

3

Exercises

- Some Stuff for the Lab



Regression

Intuition

- The regression model is a procedure that allows to estimate certain relationship that relates two or more variables with an output.

We have two types

- Linear Regression
- Non-Linear Regression



Regression

Intuition

- The regression model is a procedure that allows to estimate certain relationship that relates two or more variables with an output.

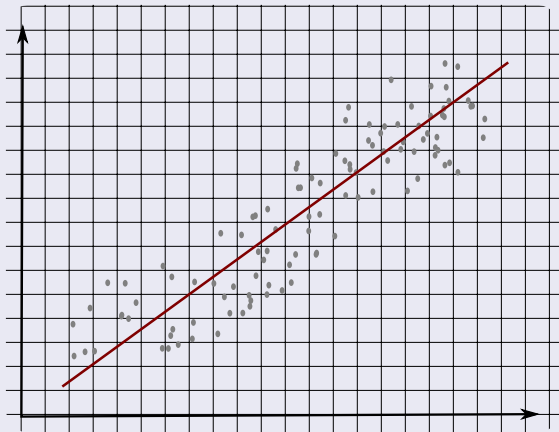
We have two types

- Linear Regression
- Non-Linear Regression



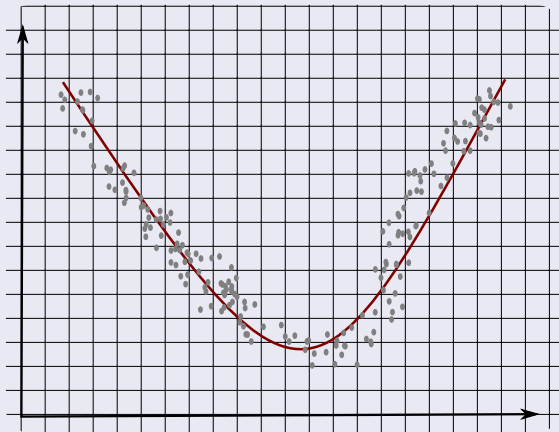
Linear Regression

We have something like



Non-Linear Regression

We have something like



As an Approximation

It is clear that in these univariate cases, we have

$$\{(x_i, y_i)\}_{i=1}^N \text{ with } x_i, y_i \in \mathbb{R}$$

Data to try to approximate by

$$\min_f \otimes_{i=1}^N g \{f(x_i) \oplus y_i\}$$

where

- \otimes, \oplus are binary operators
- g, f are functions



As an Approximation

It is clear that in these univariate cases, we have

$$\{(x_i, y_i)\}_{i=1}^N \text{ with } x_i, y_i \in \mathbb{R}$$

Data to try to approximate by

$$\min_f \otimes_{i=1}^N g \{f(x_i) \oplus y_i\}$$

Where

- \otimes, \oplus are binary operators
- g, f are functions



As an Approximation

It is clear that in these univariate cases, we have

$$\{(x_i, y_i)\}_{i=1}^N \text{ with } x_i, y_i \in \mathbb{R}$$

Data to try to approximate by

$$\min_f \otimes_{i=1}^N g \{f(x_i) \oplus y_i\}$$

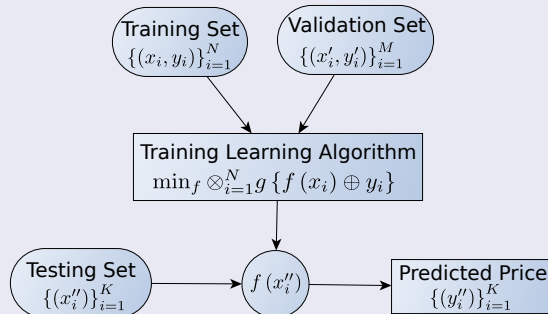
Where

- \otimes, \oplus are binary operators
- g, f are functions



Then, in Supervised Training

We have the following process (x_i, y_i)



Outline

1

Introduction

- Introduction
- Regression as approximation
- **The Simplest Functions**
- Splitting the Space
- Defining the Decision Surface
- Properties of the Hyperplane $w^T x + w_0$
- Augmenting the Vector

2

Developing a Solution

- Least Squared Error Procedure
 - The Geometry of a Two-Category Linearly-Separable Case
- The Error Idea
- The Final Error Equation
- Basic Solution
- Multidimensional Solution
- Remember in matrices of 3×3
- What Lives Where?
- Geometric Interpretation
- Solving the Labeling Issue
- Multi-Class Solution
- Issues with Least Squares!!!
 - Singularity Notes
 - Problem with Outliers
 - Problem with High Number of Dimensions
- What can be done?
 - Using Statistics to find Important Features
 - What about Numerical Stability?
 - Ridge Regression
- Observation About Eigenvalues

3

Exercises

- Some Stuff for the Lab



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:

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



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

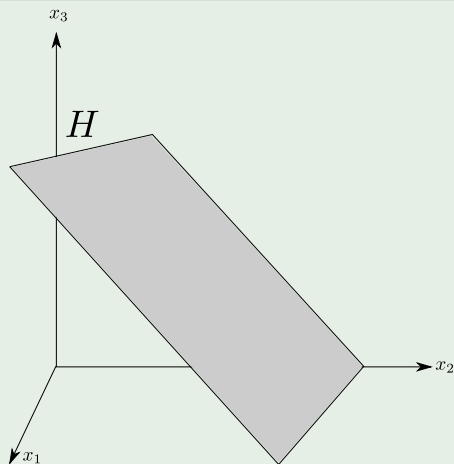
We have:

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



Example

Hyperplane in \mathbb{R}^3



Outline

1

Introduction

- Introduction
- Regression as approximation
- The Simplest Functions
- **Splitting the Space**
- Defining the Decision Surface
- Properties of the Hyperplane $w^T x + w_0$
- Augmenting the Vector

2

Developing a Solution

- Least Squared Error Procedure
 - The Geometry of a Two-Category Linearly-Separable Case
- The Error Idea
- The Final Error Equation
- Basic Solution
- Multidimensional Solution
- Remember in matrices of 3×3
- What Lives Where?
- Geometric Interpretation
- Solving the Labeling Issue
- Multi-Class Solution
- Issues with Least Squares!!!
 - Singularity Notes
 - Problem with Outliers
 - Problem with High Number of Dimensions
- What can be done?
 - Using Statistics to find Important Features
 - What about Numerical Stability?
 - Ridge Regression
- Observation About Eigenvalues

3

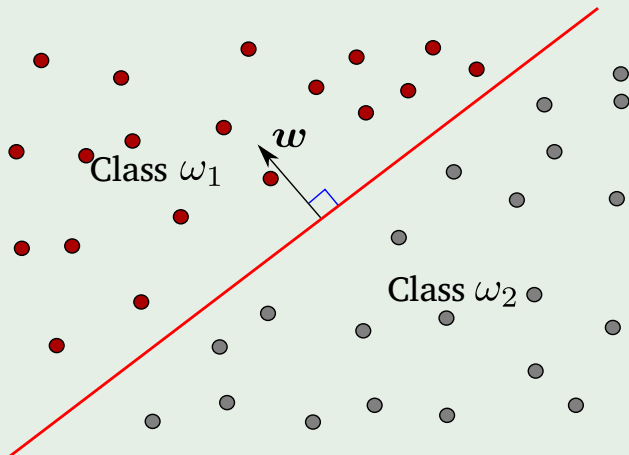
Exercises

- Some Stuff for the Lab



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

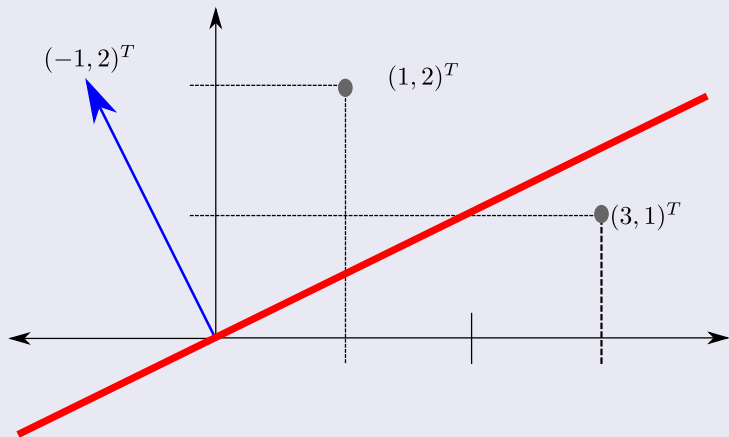


Splitting the Space?

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

Introduction

- Introduction
- Regression as approximation
- The Simplest Functions
- Splitting the Space
- **Defining the Decision Surface**
- Properties of the Hyperplane $w^T x + w_0$
- Augmenting the Vector

2

Developing a Solution

- Least Squared Error Procedure
 - The Geometry of a Two-Category Linearly-Separable Case
- The Error Idea
- The Final Error Equation
- Basic Solution
- Multidimensional Solution
- Remember in matrices of 3×3
- What Lives Where?
- Geometric Interpretation
- Solving the Labeling Issue
- Multi-Class Solution
- Issues with Least Squares!!!
 - Singularity Notes
 - Problem with Outliers
 - Problem with High Number of Dimensions
- What can be done?
 - Using Statistics to find Important Features
 - What about Numerical Stability?
 - Ridge Regression
- Observation About Eigenvalues

3

Exercises

- Some Stuff for the Lab



Cinvestav

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

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

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

Thus

$$w^T x_1 + w_0 = w^T x_2 + w_0 \quad (3)$$



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

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

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

Thus

$$w^T x_1 + w_0 = w^T x_2 + w_0 \quad (3)$$



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 \mathbf{x}_1 and \mathbf{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

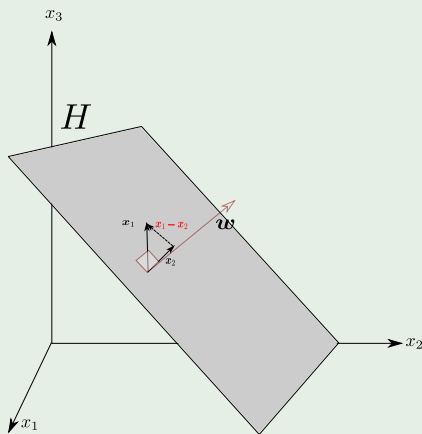
$$\mathbf{w}^T (\mathbf{x}_1 - \mathbf{x}_2) = 0 \quad (4)$$



Therefore

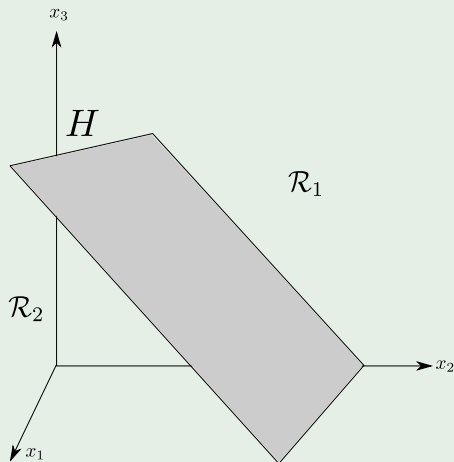
$x_1 - x_2$ lives in the hyperplane i.e. it is perpendicular to w^T

- Remark: any vector in the hyperplane is a linear combination of elements in the plane.
- **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

Introduction

- Introduction
- Regression as approximation
- The Simplest Functions
- Splitting the Space
- Defining the Decision Surface
- **Properties of the Hyperplane $w^T x + w_0$**
- Augmenting the Vector

2

Developing a Solution

- Least Squared Error Procedure
 - The Geometry of a Two-Category Linearly-Separable Case
- The Error Idea
- The Final Error Equation
- Basic Solution
- Multidimensional Solution
- Remember in matrices of 3×3
- What Lives Where?
- Geometric Interpretation
- Solving the Labeling Issue
- Multi-Class Solution
- Issues with Least Squares!!!
 - Singularity Notes
 - Problem with Outliers
 - Problem with High Number of Dimensions
- What can be done?
 - Using Statistics to find Important Features
 - What about Numerical Stability?
 - Ridge Regression
- Observation About Eigenvalues

3

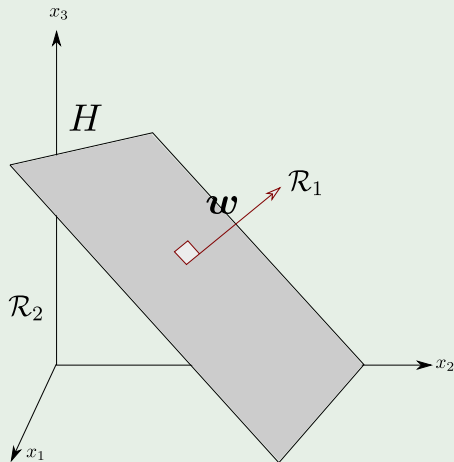
Exercises

- Some Stuff for the Lab



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

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

First, we express any \mathbf{x} as follows

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

where

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

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

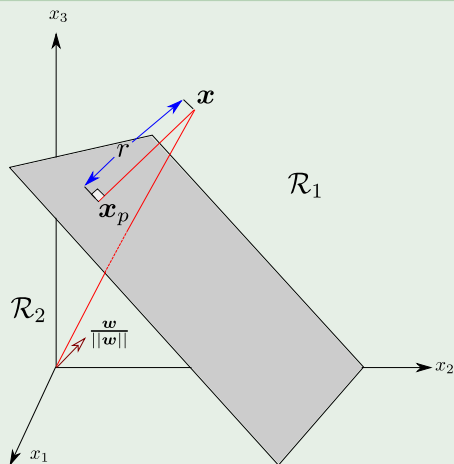
$$x = x_p + r \frac{w}{\|w\|}$$

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

Introduction

- Introduction
- Regression as approximation
- The Simplest Functions
- Splitting the Space
- Defining the Decision Surface
- Properties of the Hyperplane $w^T x + w_0$
- **Augmenting the Vector**

2

Developing a Solution

- Least Squared Error Procedure
 - The Geometry of a Two-Category Linearly-Separable Case
- The Error Idea
- The Final Error Equation
- Basic Solution
- Multidimensional Solution
- Remember in matrices of 3×3
- What Lives Where?
- Geometric Interpretation
- Solving the Labeling Issue
- Multi-Class Solution
- Issues with Least Squares!!!
 - Singularity Notes
 - Problem with Outliers
 - Problem with High Number of Dimensions
- What can be done?
 - Using Statistics to find Important Features
 - What about Numerical Stability?
 - Ridge Regression
- Observation About Eigenvalues

3

Exercises

- Some Stuff for the Lab



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

Introduction

- Introduction
- Regression as approximation
- The Simplest Functions
- Splitting the Space
- Defining the Decision Surface
- Properties of the Hyperplane $w^T x + w_0$
- Augmenting the Vector

2

Developing a Solution

● Least Squared Error Procedure

- The Geometry of a Two-Category Linearly-Separable Case
- The Error Idea
- The Final Error Equation
- Basic Solution
- Multidimensional Solution
- Remember in matrices of 3×3
- What Lives Where?
- Geometric Interpretation
- Solving the Labeling Issue
- Multi-Class Solution
- Issues with Least Squares!!!
 - Singularity Notes
 - Problem with Outliers
 - Problem with High Number of Dimensions
- What can be done?
 - Using Statistics to find Important Features
 - What about Numerical Stability?
 - Ridge Regression
- Observation About Eigenvalues

3

Exercises

- Some Stuff for the Lab



Remember

Our original function

$$\min_f \bigotimes_{i=1}^N g \{ f(x_i) \oplus y_i \}$$



Outline

1

Introduction

- Introduction
- Regression as approximation
- The Simplest Functions
- Splitting the Space
- Defining the Decision Surface
- Properties of the Hyperplane $w^T x + w_0$
- Augmenting the Vector

2

Developing a Solution

● Least Squared Error Procedure

- The Geometry of a Two-Category Linearly-Separable Case
 - The Error Idea
 - The Final Error Equation
 - Basic Solution
 - Multidimensional Solution
 - Remember in matrices of 3×3
 - What Lives Where?
 - Geometric Interpretation
 - Solving the Labeling Issue
 - Multi-Class Solution
 - Issues with Least Squares!!!
 - Singularity Notes
 - Problem with Outliers
 - Problem with High Number of Dimensions
- What can be done?
 - Using Statistics to find Important Features
 - What about Numerical Stability?
 - Ridge Regression
- Observation About Eigenvalues

3

Exercises

- Some Stuff for the Lab



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!!!

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

Introduction

- Introduction
- Regression as approximation
- The Simplest Functions
- Splitting the Space
- Defining the Decision Surface
- Properties of the Hyperplane $w^T x + w_0$
- Augmenting the Vector

2

Developing a Solution

- Least Squared Error Procedure
 - The Geometry of a Two-Category Linearly-Separable Case
- **The Error Idea**
- The Final Error Equation
- Basic Solution
- Multidimensional Solution
- Remember in matrices of 3×3
- What Lives Where?
- Geometric Interpretation
- Solving the Labeling Issue
- Multi-Class Solution
- Issues with Least Squares!!!
 - Singularity Notes
 - Problem with Outliers
 - Problem with High Number of Dimensions
- What can be done?
 - Using Statistics to find Important Features
 - What about Numerical Stability?
 - Ridge Regression
- Observation About Eigenvalues

3

Exercises

- Some Stuff for the Lab



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

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

Where the e

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

Thus, we can do the following

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



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



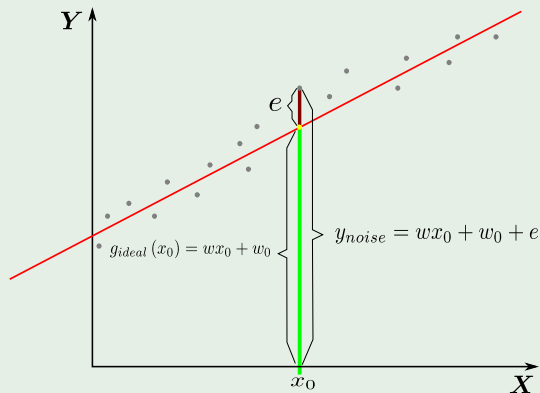
Cinvestav

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

Introduction

- Introduction
- Regression as approximation
- The Simplest Functions
- Splitting the Space
- Defining the Decision Surface
- Properties of the Hyperplane $w^T x + w_0$
- Augmenting the Vector

2

Developing a Solution

- Least Squared Error Procedure
 - The Geometry of a Two-Category Linearly-Separable Case
- The Error Idea
- **The Final Error Equation**
- Basic Solution
- Multidimensional Solution
- Remember in matrices of 3×3
- What Lives Where?
- Geometric Interpretation
- Solving the Labeling Issue
- Multi-Class Solution
- Issues with Least Squares!!!
 - Singularity Notes
 - Problem with Outliers
 - Problem with High Number of Dimensions
- What can be done?
 - Using Statistics to find Important Features
 - What about Numerical Stability?
 - Ridge Regression
- Observation About Eigenvalues

3

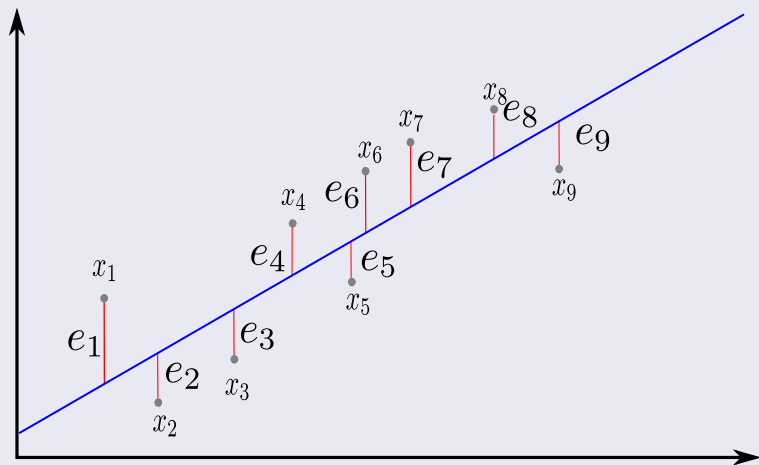
Exercises

- Some Stuff for the Lab



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 $\mathbf{x}^T = (1, \mathbf{x}')$.



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 $\mathbf{x}^T = (\mathbf{1}, \mathbf{x}')$.



Outline

1

Introduction

- Introduction
- Regression as approximation
- The Simplest Functions
- Splitting the Space
- Defining the Decision Surface
- Properties of the Hyperplane $w^T x + w_0$
- Augmenting the Vector

2

Developing a Solution

- Least Squared Error Procedure
 - The Geometry of a Two-Category Linearly-Separable Case
- The Error Idea
- The Final Error Equation
- **Basic Solution**
- Multidimensional Solution
- Remember in matrices of 3×3
- What Lives Where?
- Geometric Interpretation
- Solving the Labeling Issue
- Multi-Class Solution
- Issues with Least Squares!!!
 - Singularity Notes
 - Problem with Outliers
 - Problem with High Number of Dimensions
- What can be done?
 - Using Statistics to find Important Features
 - What about Numerical Stability?
 - Ridge Regression
- Observation About Eigenvalues

3

Exercises

- Some Stuff for the Lab



Assume that $x \in \mathbb{R}$

Then we have that the function looks like

$$f(x) = b_0 + b_1x$$

Therefore the loss function looks like

$$L(b_0, b_1, \{x_i, y_i\}_{i=1}^N) = \sum_{i=1}^N e_i^2 = \sum_{i=1}^N [y_i - b_0 - b_1x_i]^2$$



Assume that $x \in \mathbb{R}$

Then we have that the function looks like

$$f(x) = b_0 + b_1x$$

Therefore the lose function looks like

$$L(b_1, b_2, \{x_i, y_i\}_{i=1}^N) = \sum_{i=1}^N e_i^2 = \sum_{i=1}^N [y_i - b_0 - b_1x]^2$$



Then, you use simple derivatives

Then, you have derivatives with respect to b_0

$$\frac{\partial \sum_{i=1}^N e_i^2}{\partial b_0} = -2 \sum_{i=1}^N [y_i - b_0 - b_1 x] = 0$$

Derivatives with respect to b_1

$$\frac{\partial \sum_{i=1}^N e_i^2}{\partial b_1} = -2 \sum_{i=1}^N [y_i - b_0 - b_1 x] x = 0$$



Then, you use simple derivatives

Then, you have derivatives with respect to b_0

$$\frac{\partial \sum_{i=1}^N e_i^2}{\partial b_0} = -2 \sum_{i=1}^N [y_i - b_0 - b_1 x] = 0$$

Derivatives with respect to b_1

$$\frac{\partial \sum_{i=1}^N e_i^2}{\partial b_1} = -2 \sum_{i=1}^N [y_i - b_0 - b_1 x] x = 0$$



Previous equations are known as normal equations

Solving them

$$b_0 = \bar{y} - b_1 \bar{x}$$

$$b_1 = \frac{\sum_{i=1}^N [x_i - \bar{x}] [y_i - \bar{y}]}{\sum_{i=1}^N [x_i - \bar{x}]^2}$$



Outline

1

Introduction

- Introduction
- Regression as approximation
- The Simplest Functions
- Splitting the Space
- Defining the Decision Surface
- Properties of the Hyperplane $w^T x + w_0$
- Augmenting the Vector

2

Developing a Solution

- Least Squared Error Procedure
 - The Geometry of a Two-Category Linearly-Separable Case
- The Error Idea
- The Final Error Equation
- Basic Solution
- **Multidimensional Solution**
- Remember in matrices of 3×3
- What Lives Where?
- Geometric Interpretation
- Solving the Labeling Issue
- Multi-Class Solution
- Issues with Least Squares!!!
 - Singularity Notes
 - Problem with Outliers
 - Problem with High Number of Dimensions
- What can be done?
 - Using Statistics to find Important Features
 - What about Numerical Stability?
 - Ridge Regression
- Observation About Eigenvalues

3

Exercises

- Some Stuff for the Lab



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



Outline

1

Introduction

- Introduction
- Regression as approximation
- The Simplest Functions
- Splitting the Space
- Defining the Decision Surface
- Properties of the Hyperplane $w^T x + w_0$
- Augmenting the Vector

2

Developing a Solution

- Least Squared Error Procedure
 - The Geometry of a Two-Category Linearly-Separable Case
- The Error Idea
- The Final Error Equation
- Basic Solution
- Multidimensional Solution
- **Remember in matrices of 3×3**
- What Lives Where?
- Geometric Interpretation
- Solving the Labeling Issue
- Multi-Class Solution
- Issues with Least Squares!!!
 - Singularity Notes
 - Problem with Outliers
 - Problem with High Number of Dimensions
- What can be done?
 - Using Statistics to find Important Features
 - What about Numerical Stability?
 - Ridge Regression
- Observation About Eigenvalues

3

Exercises

- Some Stuff for the Lab



Some Basic Definitions

Transpose of a Matrix

$$\begin{pmatrix} a_{11} & a_{12} & a_{13} \\ a_{21} & a_{22} & a_{23} \\ a_{31} & a_{32} & a_{33} \end{pmatrix}^T = \begin{pmatrix} a_{11} & a_{21} & a_{31} \\ a_{12} & a_{22} & a_{32} \\ a_{13} & a_{23} & a_{33} \end{pmatrix}$$

$$\begin{pmatrix} a_1 \\ a_2 \\ a_3 \end{pmatrix}^T = (a_1 \ a_2 \ a_3)$$



Some Basic Definitions

Transpose of a Matrix

$$\begin{pmatrix} a_{11} & a_{12} & a_{13} \\ a_{21} & a_{22} & a_{23} \\ a_{31} & a_{32} & a_{33} \end{pmatrix}^T = \begin{pmatrix} a_{11} & a_{21} & a_{31} \\ a_{12} & a_{22} & a_{32} \\ a_{13} & a_{23} & a_{33} \end{pmatrix}$$

$$\begin{pmatrix} a_1 \\ a_2 \\ a_3 \end{pmatrix}^T = \begin{pmatrix} a_1 & a_2 & a_3 \end{pmatrix}$$



Additionally

We have

Given A and B matrices:

- $(A + B)^T = A^T + B^T$
- $(AB)^T = B^T A^T$

Given vectors x , y and a matrix A such that you can multiply them:

- $x^T A y = [x^T A y]^T = y^T A^T x$ given that the transpose of a number is the number itself.



Additionally

We have

Given A and B matrices:

- $(A + B)^T = A^T + B^T$
- $(AB)^T = B^T A^T$

Given vectors x , y and a matrix A such that you can multiply them:

- $x^T Ay = [x^T Ay]^T = y^T A^T x$ given that the transpose of a number is the number itself.



Additionally

We have

Given A and B matrices:

- $(A + B)^T = A^T + B^T$
- $(AB)^T = B^T A^T$

Given vectors x , y and a matrix A such that you can multiply them:

- $x^T A y = [x^T A y]^T = y^T A^T x$ given that the transpose of a number is the number itself.



Additionally

We have

Given A and B matrices:

- $(A + B)^T = A^T + B^T$
- $(AB)^T = B^T A^T$

Given vectors x , y and a matrix A such that you can multiply them:

- $x^T A y = \left[x^T A y \right]^T = y^T A^T x$ given that the transpose of a number is the number itself.



Additionally

We have

Given A and B matrices:

- $(A + B)^T = A^T + B^T$
- $(AB)^T = B^T A^T$

Given vectors x , y and a matrix A such that you can multiply them:

- $x^T A y = [x^T A y]^T = y^T A^T x$ given that the transpose of a number is the number itself.



Some Basic Definitions for

Derivative on Matrices

$$\frac{dA\mathbf{x}}{d\mathbf{x}} = \frac{d \begin{pmatrix} a_{11} & a_{12} & a_{13} \\ a_{21} & a_{22} & a_{23} \\ a_{31} & a_{32} & a_{33} \end{pmatrix} \begin{pmatrix} x_1 \\ x_2 \\ x_3 \end{pmatrix}}{d \begin{pmatrix} x_1 \\ x_2 \\ x_3 \end{pmatrix}}$$



Therefore

We have

$$\frac{d \begin{pmatrix} a_{11}x_1 + a_{12}x_2 + a_{13}x_3 \\ a_{21}x_1 + a_{22}x_2 + a_{23}x_3 \\ a_{31}x_1 + a_{32}x_2 + a_{33}x_3 \end{pmatrix}}{d \begin{pmatrix} x_1 \\ x_2 \\ x_3 \end{pmatrix}} = \dots$$

$$\begin{pmatrix} \frac{d(a_{11}x_1 + a_{12}x_2 + a_{13}x_3)}{dx_1} & \frac{d(a_{11}x_1 + a_{12}x_2 + a_{13}x_3)}{dx_2} & \frac{d(a_{11}x_1 + a_{12}x_2 + a_{13}x_3)}{dx_3} \\ \frac{d(a_{21}x_1 + a_{22}x_2 + a_{23}x_3)}{dx_1} & \frac{d(a_{21}x_1 + a_{22}x_2 + a_{23}x_3)}{dx_2} & \frac{d(a_{21}x_1 + a_{22}x_2 + a_{23}x_3)}{dx_3} \\ \frac{d(a_{31}x_1 + a_{32}x_2 + a_{33}x_3)}{dx_1} & \frac{d(a_{31}x_1 + a_{32}x_2 + a_{33}x_3)}{dx_2} & \frac{d(a_{31}x_1 + a_{32}x_2 + a_{33}x_3)}{dx_3} \end{pmatrix} = \dots$$

$$\begin{pmatrix} a_{11} & a_{12} & a_{13} \\ a_{21} & a_{22} & a_{23} \\ a_{31} & a_{32} & a_{33} \end{pmatrix}$$

Therefore

We have

$$\frac{d \begin{pmatrix} a_{11}x_1 + a_{12}x_2 + a_{13}x_3 \\ a_{21}x_1 + a_{22}x_2 + a_{23}x_3 \\ a_{31}x_1 + a_{32}x_2 + a_{33}x_3 \end{pmatrix}}{d \begin{pmatrix} x_1 \\ x_2 \\ x_3 \end{pmatrix}} = \dots$$

$$\left(\begin{array}{ccc} \frac{d(a_{11}x_1 + a_{12}x_2 + a_{13}x_3)}{dx_1} & \frac{d(a_{11}x_1 + a_{12}x_2 + a_{13}x_3)}{dx_2} & \frac{d(a_{11}x_1 + a_{12}x_2 + a_{13}x_3)}{dx_3} \\ \frac{d(a_{21}x_1 + a_{22}x_2 + a_{23}x_3)}{dx_1} & \frac{d(a_{21}x_1 + a_{22}x_2 + a_{23}x_3)}{dx_2} & \frac{d(a_{21}x_1 + a_{22}x_2 + a_{23}x_3)}{dx_3} \\ \frac{d(a_{31}x_1 + a_{32}x_2 + a_{33}x_3)}{dx_1} & \frac{d(a_{31}x_1 + a_{32}x_2 + a_{33}x_3)}{dx_2} & \frac{d(a_{31}x_1 + a_{32}x_2 + a_{33}x_3)}{dx_3} \end{array} \right) = \dots$$

$$\begin{pmatrix} a_{11} & a_{12} & a_{13} \\ a_{21} & a_{22} & a_{23} \\ a_{31} & a_{32} & a_{33} \end{pmatrix}$$

Therefore

We have

$$\frac{d \begin{pmatrix} a_{11}x_1 + a_{12}x_2 + a_{13}x_3 \\ a_{21}x_1 + a_{22}x_2 + a_{23}x_3 \\ a_{31}x_1 + a_{32}x_2 + a_{33}x_3 \end{pmatrix}}{d \begin{pmatrix} x_1 \\ x_2 \\ x_3 \end{pmatrix}} = \dots$$

$$\left(\begin{array}{ccc} \frac{d(a_{11}x_1 + a_{12}x_2 + a_{13}x_3)}{dx_1} & \frac{d(a_{11}x_1 + a_{12}x_2 + a_{13}x_3)}{dx_2} & \frac{d(a_{11}x_1 + a_{12}x_2 + a_{13}x_3)}{dx_3} \\ \frac{d(a_{21}x_1 + a_{22}x_2 + a_{23}x_3)}{dx_1} & \frac{d(a_{21}x_1 + a_{22}x_2 + a_{23}x_3)}{dx_2} & \frac{d(a_{21}x_1 + a_{22}x_2 + a_{23}x_3)}{dx_3} \\ \frac{d(a_{31}x_1 + a_{32}x_2 + a_{33}x_3)}{dx_1} & \frac{d(a_{31}x_1 + a_{32}x_2 + a_{33}x_3)}{dx_2} & \frac{d(a_{31}x_1 + a_{32}x_2 + a_{33}x_3)}{dx_3} \end{array} \right) = \dots$$

$$\begin{pmatrix} a_{11} & a_{12} & a_{13} \\ a_{21} & a_{22} & a_{23} \\ a_{31} & a_{32} & a_{33} \end{pmatrix}$$

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.

Wait. How we get the discriminant function?

Any Ideas?



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.

Thus, How we get the discriminant function?

Any Ideas?



The Final Discriminant Function

Very Simple!!!

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



Pseudo-inverse of a Matrix

Definition

Suppose that $X \in \mathbb{R}^{m \times n}$ and $\text{rank}(X) = m$. We call the matrix

$$X^+ = (X^T X)^{-1} X^T$$

the pseudo inverse of X .

Pseudo-inverse of a Matrix

Definition

Suppose that $X \in \mathbb{R}^{m \times n}$ and $\text{rank}(X) = m$. We call the matrix

$$\mathbf{X}^+ = (\mathbf{X}^T \mathbf{X})^{-1} \mathbf{X}^T$$

the pseudo inverse of X .

Reason

X^+ inverts X on its image

Pseudo-inverse of a Matrix

Definition

Suppose that $X \in \mathbb{R}^{m \times n}$ and $\text{rank}(X) = m$. We call the matrix

$$X^+ = (X^T X)^{-1} X^T$$

the pseudo inverse of X .

Reason

X^+ inverts X on its image

What?

- First a definition
 - ▶ If $w \in \text{image}(X)$, then there is some $v \in \mathbb{R}^n$ such that $w = Xv$.
- Hence, $X^+w = X^+Xv = (X^T X)^{-1} X^T Xv = v$

Pseudo-inverse of a Matrix

Definition

Suppose that $X \in \mathbb{R}^{m \times n}$ and $\text{rank}(X) = m$. We call the matrix

$$X^+ = (X^T X)^{-1} X^T$$

the pseudo inverse of X .

Reason

X^+ inverts X on its image

What?

- First a definition
 - ▶ If $w \in \text{image}(X)$, then there is some $v \in \mathbb{R}^n$ such that $w = Xv$.
- Hence, $X^+w = X^+Xv = (X^T X)^{-1} X^T Xv = v$

Pseudo-inverse of a Matrix

Definition

Suppose that $X \in \mathbb{R}^{m \times n}$ and $\text{rank}(X) = m$. We call the matrix

$$X^+ = (X^T X)^{-1} X^T$$

the pseudo inverse of X .

Reason

X^+ inverts X on its image

What?

- First a definition
 - ▶ If $w \in \text{image}(X)$, then there is some $v \in \mathbb{R}^n$ such that $w = Xv$.
- Hence, $X^+w = X^+Xv = (X^T X)^{-1} X^T Xv = v$

Pseudo-inverse of a Matrix

Definition

Suppose that $X \in \mathbb{R}^{m \times n}$ and $\text{rank}(X) = m$. We call the matrix

$$X^+ = (X^T X)^{-1} X^T$$

the pseudo inverse of X .

Reason

X^+ inverts X on its image

What?

- First a definition
 - ▶ If $w \in \text{image}(X)$, then there is some $v \in \mathbb{R}^n$ such that $w = Xv$.

• Hence, $X^+w = X^+Xv = (X^T X)^{-1} X^T Xv = v$

Pseudo-inverse of a Matrix

Definition

Suppose that $X \in \mathbb{R}^{m \times n}$ and $\text{rank}(X) = m$. We call the matrix

$$X^+ = (X^T X)^{-1} X^T$$

the pseudo inverse of X .

Reason

X^+ inverts X on its image

What?

- First a definition
 - ▶ If $w \in \text{image}(X)$, then there is some $v \in \mathbb{R}^n$ such that $w = Xv$.
- Hence, $X^+w = X^+Xv = (X^T X)^{-1} X^T Xv = v$

Outline

1

Introduction

- Introduction
- Regression as approximation
- The Simplest Functions
- Splitting the Space
- Defining the Decision Surface
- Properties of the Hyperplane $w^T x + w_0$
- Augmenting the Vector

2

Developing a Solution

- Least Squared Error Procedure
 - The Geometry of a Two-Category Linearly-Separable Case
- The Error Idea
- The Final Error Equation
- Basic Solution
- Multidimensional Solution
- Remember in matrices of 3×3
- **What Lives Where?**
 - Geometric Interpretation
 - Solving the Labeling Issue
 - Multi-Class Solution
 - Issues with Least Squares!!!
 - Singularity Notes
 - Problem with Outliers
 - Problem with High Number of Dimensions
- What can be done?
 - Using Statistics to find Important Features
 - What about Numerical Stability?
 - Ridge Regression
- Observation About Eigenvalues

3

Exercises

- Some Stuff for the Lab



We have that

The Data Matrix

$$\mathbf{X} \in \mathbb{R}^{N \times (d+1)}$$

$$x_i \in \mathbb{R}^d$$



We have that

The Data Matrix

$$\mathbf{X} \in \mathbb{R}^{N \times (d+1)}$$

$$x_i \in \mathbb{R}^d$$



We have that

The Data Matrix

$$\mathbf{X} \in \mathbb{R}^{N \times (d+1)}$$

$$\mathbf{x}_i \in \mathbb{R}^d$$



The projected elements by the matrix

Definition *Image* (\mathbf{X})

- The column space of a matrix \mathbf{X} is the span (set of all possible linear combinations) of its column vectors.

$$\text{Image}(\mathbf{X}) = \text{span} \{ \mathbf{X}_1^{\text{col}}, \dots, \mathbf{X}_{d+1}^{\text{col}} \}$$

- ▶ In other words, the image of a matrix \mathbf{X} is all the vectors $\mathbf{X}\mathbf{v} \in \mathbb{R}^N$ with $\mathbf{v} \in \mathbb{R}^{d+1}$



The Data Samples

The Data Samples

$$\mathbf{x}_i \in \mathbb{R}^d$$



Additionally, we have that

The Weight Vector w

$$w \in \mathbb{R}^{d+1}$$

What about the column space of X and the ideal input vector y

$$X_i^{col}, y \in \mathbb{R}^N$$



Additionally, we have that

The Weight Vector w

$$w \in \mathbb{R}^{d+1}$$

What about the column space of X and the ideal input vector y

$$X_i^{col}, y \in \mathbb{R}^N$$



We can now see where \mathbf{y} is being projected

Basically \mathbf{y} , the list of real inputs is being projected into

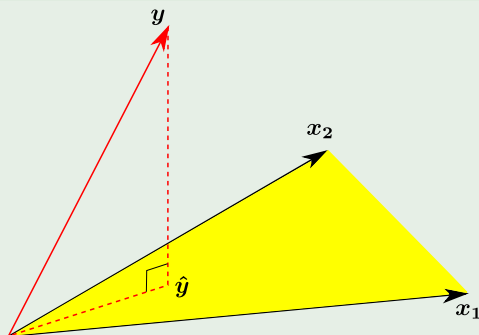
$$\text{span} \{ \mathbf{x}_1, \mathbf{x}_2, \dots, \mathbf{x}_N \} \quad (22)$$

- by function $\hat{\mathbf{y}} = \mathbf{X} (\mathbf{X}^T \mathbf{X})^{-1} \mathbf{X}^T \mathbf{y}$.



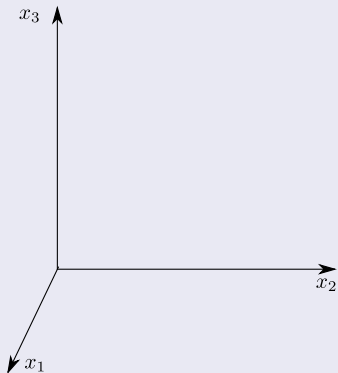
Geometrically

Given a y , you obtain a projected \hat{y} through the projection function $X (X^T X)^{-1} X^T$



Why? Assume that you are in \mathbb{R}^3

Something like



Simple but complex

A simple question

- What are the projections of $b = (2, 3, 4)$ onto the z axis and the xy plane?
- Can we use matrices to talk about these projections?

First

We must have a projection matrix P with the following property:

$$P^2 = P$$

Why?

Ideas?



Cinvestav

Simple but complex

A simple question

- What are the projections of $b = (2, 3, 4)$ onto the z axis and the xy plane?
- Can we use matrices to talk about these projections?

First

We must have a projection matrix P with the following property:

$$P^2 = P$$

Ideas?



Cinvestav

Simple but complex

A simple question

- What are the projections of $b = (2, 3, 4)$ onto the z axis and the xy plane?
- Can we use matrices to talk about these projections?

First

We must have a projection matrix P with the following property:

$$P^2 = P$$

Why?

Ideas?



Then, the Projection Pb

First

When b is projected onto a line, its projection p is the part of b along that line.

Second

When b is projected onto a plane, its projection p is the part of the plane.



Then, the Projection $P\mathbf{b}$

First

When \mathbf{b} is projected onto a line, its projection \mathbf{p} is the part of \mathbf{b} along that line.

Second

When \mathbf{b} is projected onto a plane, its projection \mathbf{p} is the part of the plane.



In our case

The Projection Matrices for the coordinate systems

$$P_1 = \begin{pmatrix} 0 & 0 & 0 \\ 0 & 0 & 0 \\ 0 & 0 & 1 \end{pmatrix}, P_2 = \begin{pmatrix} 0 & 0 & 0 \\ 0 & 1 & 0 \\ 0 & 0 & 0 \end{pmatrix}, P_3 = \begin{pmatrix} 1 & 0 & 0 \\ 0 & 0 & 0 \\ 0 & 0 & 0 \end{pmatrix}$$



Example

We have the following vector $\mathbf{b} = (2, 3, 4)^T$

Onto the z axis:

$$P_1 \mathbf{b} = \begin{pmatrix} 0 & 0 & 0 \\ 0 & 0 & 0 \\ 0 & 0 & 1 \end{pmatrix} \begin{pmatrix} 2 \\ 3 \\ 4 \end{pmatrix} = \begin{pmatrix} 0 \\ 0 \\ 4 \end{pmatrix}$$

What about the plane xy ?

Any idea?



Example

We have the following vector $\mathbf{b} = (2, 3, 4)^T$

Onto the z axis:

$$P_1 \mathbf{b} = \begin{pmatrix} 0 & 0 & 0 \\ 0 & 0 & 0 \\ 0 & 0 & 1 \end{pmatrix} \begin{pmatrix} 2 \\ 3 \\ 4 \end{pmatrix} = \begin{pmatrix} 0 \\ 0 \\ 4 \end{pmatrix}$$

What about the plane xy

Any idea?



We have something more complex

Something Notable

$$P_4 = \begin{pmatrix} 1 & 0 & 0 \\ 0 & 1 & 0 \\ 0 & 0 & 0 \end{pmatrix}$$

Then

$$P_4 b = \begin{pmatrix} 1 & 0 & 0 \\ 0 & 1 & 0 \\ 0 & 0 & 0 \end{pmatrix} \begin{pmatrix} 2 \\ 3 \\ 4 \end{pmatrix} = \begin{pmatrix} 2 \\ 3 \\ 0 \end{pmatrix}$$



We have something more complex

Something Notable

$$P_4 = \begin{pmatrix} 1 & 0 & 0 \\ 0 & 1 & 0 \\ 0 & 0 & 0 \end{pmatrix}$$

Then

$$P_4 \mathbf{b} = \begin{pmatrix} 1 & 0 & 0 \\ 0 & 1 & 0 \\ 0 & 0 & 0 \end{pmatrix} \begin{pmatrix} 2 \\ 3 \\ 4 \end{pmatrix} = \begin{pmatrix} 2 \\ 3 \\ 0 \end{pmatrix}$$



Assume the following

We have that

$\mathbf{a}_1, \mathbf{a}_2, \dots, \mathbf{a}_n$ in \mathbb{R}^m .

Assume they are linearly independent.

They span a subspace, we want projections into the subspace.

We want to project b into such subspace.

How do we do it?



Assume the following

We have that

$\mathbf{a}_1, \mathbf{a}_2, \dots, \mathbf{a}_n$ in \mathbb{R}^m .

Assume they are linearly independent

They span a subspace, we want projections into the subspace

We want to project b into such subspace

How do we do it?



Assume the following

We have that

$\mathbf{a}_1, \mathbf{a}_2, \dots, \mathbf{a}_n$ in \mathbb{R}^m .

Assume they are linearly independent

They span a subspace, we want projections into the subspace

We want to project \mathbf{b} into such subspace

How do we do it?



This is the important part

Problem

Find the combination $\mathbf{p} = x_1\mathbf{a}_1 + x_2\mathbf{a}_2 + \cdots + x_n\mathbf{a}_n$ closest to vector \mathbf{b} .

Something Notable

With $n = 1$ (only one vector \mathbf{a}_1) this projection onto a line.

This line is the column space of \mathbf{A} .

Basically the columns are spanned by a single column.



This is the important part

Problem

Find the combination $\mathbf{p} = x_1\mathbf{a}_1 + x_2\mathbf{a}_2 + \cdots + x_n\mathbf{a}_n$ closest to vector \mathbf{b} .

Something Notable

With $n = 1$ (only one vector \mathbf{a}_1) this projection onto a line.

This line is the column space of \mathbf{A} .

Basically the columns are spanned by a single column.



This is the important part

Problem

Find the combination $\mathbf{p} = x_1\mathbf{a}_1 + x_2\mathbf{a}_2 + \cdots + x_n\mathbf{a}_n$ closest to vector \mathbf{b} .

Something Notable

With $n = 1$ (only one vector \mathbf{a}_1) this projection onto a line.

This line is the column space of A

Basically the columns are spanned by a single column.



In General

The matrix has n columns $\mathbf{a}_1, \mathbf{a}_2, \dots, \mathbf{a}_n$

The combinations in \mathbb{R}^m are vectors $A\mathbf{x}$ in the column space

We are looking for the particular combination

The nearest to the original \mathbf{b}

$$\mathbf{p} = A\hat{\mathbf{x}}$$



In General

The matrix has n columns $\mathbf{a}_1, \mathbf{a}_2, \dots, \mathbf{a}_n$

The combinations in \mathbb{R}^m are vectors $A\mathbf{x}$ in the column space

We are looking for the particular combination

The nearest to the original \mathbf{b}

$$\mathbf{p} = A\hat{\mathbf{x}}$$



First

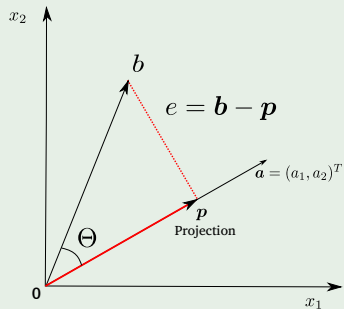
We look at the simplest case

The projection into a line...



With a little of Geometry

We have the following



Therefore

Using the fact that the projection is equal to

$$p = xa$$

Then, the error is equal to

$$e = b - xa$$

We have that $a \cdot e = 0$

$$a \cdot e = a \cdot (b - xa) = a \cdot b - xa \cdot a = 0$$



Therefore

Using the fact that the projection is equal to

$$p = xa$$

Then, the error is equal to

$$e = b - xa$$

We have that $a \cdot e = 0$

$$a \cdot e = a \cdot (b - xa) = a \cdot b - xa \cdot a = 0$$



Therefore

Using the fact that the projection is equal to

$$p = xa$$

Then, the error is equal to

$$e = b - xa$$

We have that $a \cdot e = 0$

$$a \cdot e = a \cdot (b - xa) = a \cdot b - xa \cdot a = 0$$



Therefore

We have that

$$x = \frac{\mathbf{a} \cdot \mathbf{b}}{\mathbf{a} \cdot \mathbf{a}} = \frac{\mathbf{a}^T \mathbf{b}}{\mathbf{a}^T \mathbf{a}}$$

Or something quite simple

$$p = \frac{\mathbf{a}^T \mathbf{b}}{\mathbf{a}^T \mathbf{a}} \mathbf{a}$$



Therefore

We have that

$$x = \frac{\mathbf{a} \cdot \mathbf{b}}{\mathbf{a} \cdot \mathbf{a}} = \frac{\mathbf{a}^T \mathbf{b}}{\mathbf{a}^T \mathbf{a}}$$

Or something quite simple

$$\mathbf{p} = \frac{\mathbf{a}^T \mathbf{b}}{\mathbf{a}^T \mathbf{a}} \mathbf{a}$$



By the Law of Cosines

Something Notable

$$\|\mathbf{a} - \mathbf{b}\|^2 = \|\mathbf{a}\|^2 + \|\mathbf{b}\|^2 - 2 \|\mathbf{a}\| \|\mathbf{b}\| \cos \Theta$$



We have

The following product

$$\mathbf{a} \cdot \mathbf{a} - 2\mathbf{a} \cdot \mathbf{b} + \mathbf{b} \cdot \mathbf{b} = \|\mathbf{a}\|^2 + \|\mathbf{b}\|^2 - 2\|\mathbf{a}\| \|\mathbf{b}\| \cos \Theta$$

Then

$$\mathbf{a} \cdot \mathbf{b} = \|\mathbf{a}\| \|\mathbf{b}\| \cos \Theta$$



We have

The following product

$$\mathbf{a} \cdot \mathbf{a} - 2\mathbf{a} \cdot \mathbf{b} + \mathbf{b} \cdot \mathbf{b} = \|\mathbf{a}\|^2 + \|\mathbf{b}\|^2 - 2\|\mathbf{a}\| \|\mathbf{b}\| \cos \Theta$$

Then

$$\mathbf{a} \cdot \mathbf{b} = \|\mathbf{a}\| \|\mathbf{b}\| \cos \Theta$$



With Length

Using the Norm

$$\|p\| = \left| \frac{\mathbf{a}^T \mathbf{b}}{\mathbf{a}^T \mathbf{a}} \right| \|\mathbf{a}\| = \left| \frac{\|\mathbf{a}\| \|\mathbf{b}\| \cos \Theta}{\|\mathbf{a}\|^2} \right| \|\mathbf{a}\| = \|\mathbf{b}\| |\cos \Theta|$$



Example

Project

$$\mathbf{b} = \begin{pmatrix} 1 \\ 1 \\ 1 \end{pmatrix} \text{ onto } \mathbf{a} = \begin{pmatrix} 1 \\ 2 \\ 2 \end{pmatrix}$$

Find

$$p = ra$$



Example

Project

$$\mathbf{b} = \begin{pmatrix} 1 \\ 1 \\ 1 \end{pmatrix} \text{ onto } \mathbf{a} = \begin{pmatrix} 1 \\ 2 \\ 2 \end{pmatrix}$$

Find

$$\mathbf{p} = x\mathbf{a}$$



What about the Projection Matrix in general

We have

$$\mathbf{p} = \mathbf{a}x = \frac{\mathbf{a}\mathbf{a}^T\mathbf{b}}{\mathbf{a}^T\mathbf{a}} = P\mathbf{b}$$

Then

$$P = \frac{\mathbf{a}\mathbf{a}^T}{\mathbf{a}^T\mathbf{a}}$$



What about the Projection Matrix in general

We have

$$\mathbf{p} = \mathbf{a}x = \frac{\mathbf{a}\mathbf{a}^T\mathbf{b}}{\mathbf{a}^T\mathbf{a}} = P\mathbf{b}$$

Then

$$P = \frac{\mathbf{a}\mathbf{a}^T}{\mathbf{a}^T\mathbf{a}}$$



Example

Find the projection matrix for

$$\mathbf{b} = \begin{pmatrix} 1 \\ 1 \\ 1 \end{pmatrix} \text{ onto } \mathbf{a} = \begin{pmatrix} 1 \\ 2 \\ 2 \end{pmatrix}$$



What about the general case?

We have that

Find the combination $\mathbf{p} = x_1\mathbf{a}_1 + x_2\mathbf{a}_2 + \cdots + x_n\mathbf{a}_n$ closest to vector \mathbf{b} .

Now, you need a vector

Find the vector \mathbf{x} , find the projection $\mathbf{p} = A\mathbf{x}$, find the matrix P .

Again, the error is perpendicular to the space

$$\mathbf{e} = \mathbf{b} - A\mathbf{x}$$



What about the general case?

We have that

Find the combination $\mathbf{p} = x_1\mathbf{a}_1 + x_2\mathbf{a}_2 + \cdots + x_n\mathbf{a}_n$ closest to vector \mathbf{b} .

Now you need a vector

Find the vector \mathbf{x} , find the projection $\mathbf{p} = A\mathbf{x}$, find the matrix P .

Again, the error is perpendicular to the space

$$\mathbf{e} = \mathbf{b} - A\mathbf{x}$$



What about the general case?

We have that

Find the combination $\mathbf{p} = x_1\mathbf{a}_1 + x_2\mathbf{a}_2 + \cdots + x_n\mathbf{a}_n$ closest to vector \mathbf{b} .

Now you need a vector

Find the vector \mathbf{x} , find the projection $\mathbf{p} = A\mathbf{x}$, find the matrix P .

Again, the error is perpendicular to the space

$$\mathbf{e} = \mathbf{b} - A\mathbf{x}$$



Therefore

The error $e = \mathbf{b} - A\mathbf{x}$

$$\mathbf{a}_1^T (\mathbf{b} - A\mathbf{x}) = 0$$

$$\vdots$$

$$\mathbf{a}_n^T (\mathbf{b} - A\mathbf{x}) = 0$$

Or

$$\begin{bmatrix} \mathbf{a}_1^T \\ \vdots \\ \mathbf{a}_n^T \end{bmatrix} [\mathbf{b} - A\mathbf{x}] = 0$$



Therefore

The error $e = \mathbf{b} - A\mathbf{x}$

$$\mathbf{a}_1^T (\mathbf{b} - A\mathbf{x}) = 0$$

$$\vdots$$

$$\mathbf{a}_n^T (\mathbf{b} - A\mathbf{x}) = 0$$

Or

$$\begin{bmatrix} \mathbf{a}_1^T \\ \vdots \\ \mathbf{a}_n^T \end{bmatrix} [\mathbf{b} - A\mathbf{x}] = 0$$



Therefore

The Matrix with those rows is A^T

$$A^T (\mathbf{b} - A\mathbf{x}) = 0$$

Therefore

$$A^T \mathbf{b} - A^T A \mathbf{x} = 0$$

Or the most know form

$$\mathbf{x} = (A^T A)^{-1} A^T \mathbf{b}$$



Therefore

The Matrix with those rows is A^T

$$A^T (\mathbf{b} - A\mathbf{x}) = 0$$

Therefore

$$A^T \mathbf{b} - A^T A \mathbf{x} = 0$$

Or, in the most know form

$$\mathbf{x} = (A^T A)^{-1} A^T \mathbf{b}$$



Therefore

The Matrix with those rows is A^T

$$A^T (\mathbf{b} - A\mathbf{x}) = 0$$

Therefore

$$A^T \mathbf{b} - A^T A \mathbf{x} = 0$$

Or the most know form

$$\mathbf{x} = (A^T A)^{-1} A^T \mathbf{b}$$



Therefore

The Projection is

$$\mathbf{p} = A\mathbf{x} = A \left(A^T A \right)^{-1} A^T \mathbf{b}$$

Therefore

$$P = A \left(A^T A \right)^{-1} A^T$$



Therefore

The Projection is

$$\mathbf{p} = A\mathbf{x} = A \left(A^T A \right)^{-1} A^T \mathbf{b}$$

Therefore

$$P = A \left(A^T A \right)^{-1} A^T$$



The key step was $A^T [\mathbf{b} - A\mathbf{x}] = 0$

Linear algebra gives this "normal equation"

- 1 Our subspace is the column space of A .
- 2 The error vector $\mathbf{b} - A\mathbf{x}$ is perpendicular to that column space.
- 3 Therefore $\mathbf{b} - A\mathbf{x}$ is in the nullspace of A^T



When A has independent columns, $A^T A$ is invertible

Theorem

$A^T A$ is invertible if and only if A has linearly independent columns.



Proof

Consider the following

$$A^T A \mathbf{x} = 0$$

Here, \mathbf{x} is in the null space of A^T .

- Remember the column space and null space of A^T are orthogonal complements.

And, \mathbf{x} is an element in the column space of A .

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



Proof

Consider the following

$$A^T Ax = 0$$

Here, Ax is in the null space of A^T

- Remember the column space and null space of A^T are orthogonal complements.

And, x is an element in the column space of A

$$Ax = 0$$



Proof

Consider the following

$$A^T Ax = 0$$

Here, Ax is in the null space of A^T

- Remember the column space and null space of A^T are orthogonal complements.

And Ax an element in the column space of A

$$Ax = 0$$



Proof

If A has linearly independent columns

$$Ax = 0 \implies x = 0$$

Then, the null space

$$\text{Null}(A^T A) = \{0\}$$

i.e. A is full rank

- Then, $A^T A$ is invertible...



Proof

If A has linearly independent columns

$$Ax = 0 \implies x = 0$$

Then, the null space

$$\text{Null}(A^T A) = \{0\}$$

$\Rightarrow A^T A$ is full rank

- Then, $A^T A$ is invertible...



Proof

If A has linearly independent columns

$$Ax = 0 \implies x = 0$$

Then, the null space

$$\text{Null}(A^T A) = \{0\}$$

i.e. $A^T A$ is full rank

- Then, $A^T A$ is invertible...



Finally

Theorem

- When A has independent columns, $A^T A$ is square, symmetric and invertible.



Outline

1

Introduction

- Introduction
- Regression as approximation
- The Simplest Functions
- Splitting the Space
- Defining the Decision Surface
- Properties of the Hyperplane $w^T x + w_0$
- Augmenting the Vector

2

Developing a Solution

- Least Squared Error Procedure
 - The Geometry of a Two-Category Linearly-Separable Case
- The Error Idea
- The Final Error Equation
- Basic Solution
- Multidimensional Solution
- Remember in matrices of 3×3
- What Lives Where?
- **Geometric Interpretation**
- Solving the Labeling Issue
- Multi-Class Solution
- Issues with Least Squares!!!
 - Singularity Notes
 - Problem with Outliers
 - Problem with High Number of Dimensions
- What can be done?
 - Using Statistics to find Important Features
 - What about Numerical Stability?
 - Ridge Regression
- Observation About Eigenvalues

3

Exercises

- Some Stuff for the Lab



Geometric Interpretation

We have

The image of the mapping:

$$h : \mathbf{w} \mapsto \mathbf{X}\mathbf{w}$$

$$h : \mathbb{R}^{d+1} \mapsto \mathbb{R}^N$$

is a linear subspace of \mathbb{R}^N .



What about w ?

w moves through all points in \mathbb{R}^{d+1} when being generated

- Thus, the function value $h(w) = Xw$ can move through all points in the image space:

$$\text{image}(X) = \text{span} \left\{ X_1^{\text{col}}, X_2^{\text{col}}, \dots, X_{d+1}^{\text{col}} \right\}$$

Additionally, each w defines one point in

$$h(w) = Xw = \sum_{i=1}^{d+1} w_i X_i^{\text{col}}.$$



What about w ?

w moves through all points in \mathbb{R}^{d+1} when being generated

- Thus, the function value $h(w) = Xw$ can move through all points in the image space:

$$\text{image}(X) = \text{span} \left\{ X_1^{\text{col}}, X_2^{\text{col}}, \dots, X_{d+1}^{\text{col}} \right\}$$

Additionally, each w defines one point in

$$\text{span} \left\{ X_1^{\text{col}}, X_2^{\text{col}}, \dots, X_{d+1}^{\text{col}} \right\} \subseteq \mathbb{R}^N$$

$$h(w) = Xw = \sum_{i=1}^{d+1} w_i X_i^{\text{col}}.$$



What about the optimality of w ?

We have a composition of functions that are convex

$$f(w) = w^T x$$

$$g(t) = (y - t)$$

$$h(e) = \sum_{i=1}^n e^2$$

- Making the Least Squared Error a Convex function with a single minimum!!!

The derivative method produces w

- Such that \hat{w} minimizes the distance $d(y, \text{image}(X))$.

What about the optimality of w ?

We have a composition of functions that are convex

$$f(w) = w^T x$$

$$g(t) = (y - t)$$

$$h(e) = \sum_{i=1}^n e^2$$

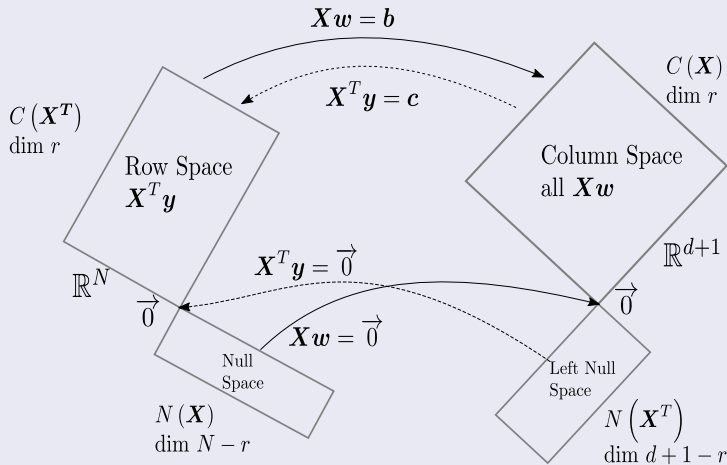
- Making the Least Squared Error a Convex function with a single minimum!!!

The derivative method produces a \hat{w}

- Such that \hat{w} minimizes the distance $d(y, \text{image}(X))$.

This comes from the following representation

Given a matrix X (“Linear Algebra and Its Applications” by Gilbert Strang)



Outline

1

Introduction

- Introduction
- Regression as approximation
- The Simplest Functions
- Splitting the Space
- Defining the Decision Surface
- Properties of the Hyperplane $w^T x + w_0$
- Augmenting the Vector

2

Developing a Solution

- Least Squared Error Procedure
 - The Geometry of a Two-Category Linearly-Separable Case
- The Error Idea
- The Final Error Equation
- Basic Solution
- Multidimensional Solution
- Remember in matrices of 3×3
- What Lives Where?
- Geometric Interpretation
- **Solving the Labeling Issue**
- Multi-Class Solution
- Issues with Least Squares!!!
 - Singularity Notes
 - Problem with Outliers
 - Problem with High Number of Dimensions
- What can be done?
 - Using Statistics to find Important Features
 - What about Numerical Stability?
 - Ridge Regression
- Observation About Eigenvalues

3

Exercises

- Some Stuff for the Lab



This Resolve Our Problem

With the Labels being chosen at the beginning

Question? Did you noticed the following?

We assume a similar number of elements in both classes

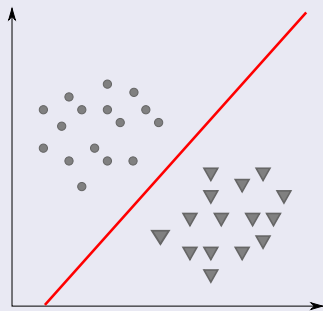


This Resolve Our Problem

With the Labels being chosen at the beginning

Question? Did you noticed the following?

We assume a similar number of elements in both classes



Outline

1

Introduction

- Introduction
- Regression as approximation
- The Simplest Functions
- Splitting the Space
- Defining the Decision Surface
- Properties of the Hyperplane $w^T x + w_0$
- Augmenting the Vector

2

Developing a Solution

- Least Squared Error Procedure
 - The Geometry of a Two-Category Linearly-Separable Case
- The Error Idea
- The Final Error Equation
- Basic Solution
- Multidimensional Solution
- Remember in matrices of 3×3
- What Lives Where?
- Geometric Interpretation
- Solving the Labeling Issue
- **Multi-Class Solution**
- Issues with Least Squares!!!
 - Singularity Notes
 - Problem with Outliers
 - Problem with High Number of Dimensions
- What can be done?
 - Using Statistics to find Important Features
 - What about Numerical Stability?
 - Ridge Regression
- Observation About Eigenvalues

3

Exercises

- Some Stuff for the Lab



Multi-Class Solution

What to do?

- 1 We might reduce the problem to $c - 1$ two-class problems.
- 2 We might use $\frac{c(c-1)}{2}$ linear discriminants, one for every pair of classes.



Multi-Class Solution

What to do?

- 1 We might reduce the problem to $c - 1$ two-class problems.
- 2 We might use $\frac{c(c-1)}{2}$ linear discriminants, one for every pair of classes.

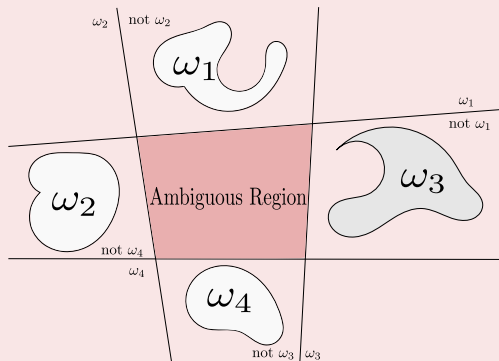


Multi-Class Solution

What to do?

- 1 We might reduce the problem to $c - 1$ two-class problems.
- 2 We might use $\frac{c(c-1)}{2}$ linear discriminants, one for every pair of classes.

However



What to Do?

Define c linear discriminant functions

$$g_i(\mathbf{x}) = \mathbf{w}^T \mathbf{x} + w_{i0} \text{ for } i = 1, \dots, c \quad (23)$$

This is known as a linear machine

Rule: if $g_k(\mathbf{x}) > g_j(\mathbf{x})$ for all $j \neq k \implies \mathbf{x} \in \omega_k$



What to Do?

Define c linear discriminant functions

$$g_i(\mathbf{x}) = \mathbf{w}^T \mathbf{x} + w_{i0} \text{ for } i = 1, \dots, c \quad (23)$$

This is known as a **linear machine**

Rule: if $g_k(\mathbf{x}) > g_j(\mathbf{x})$ for all $j \neq k \implies \mathbf{x} \in \omega_k$

Five Properties (It can be proved!!!)

- Decision Regions are Singly Connected.



What to Do?

Define c linear discriminant functions

$$g_i(\mathbf{x}) = \mathbf{w}^T \mathbf{x} + w_{i0} \text{ for } i = 1, \dots, c \quad (23)$$

This is known as a **linear machine**

Rule: if $g_k(\mathbf{x}) > g_j(\mathbf{x})$ for all $j \neq k \implies \mathbf{x} \in \omega_k$

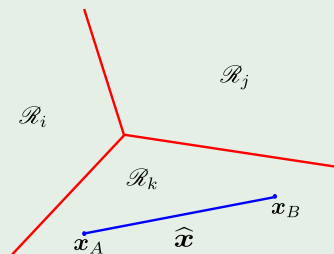
Nice Properties (It can be proved!!!)

- 1 Decision Regions are Singly Connected.
- 2 Decision Regions are Convex.



Proof of Properties

Proof



Actually quite simple

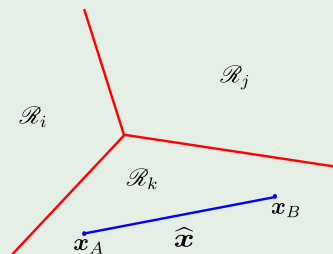
Given

$$y = \lambda x_A + (1 - \lambda) x_B$$

with $\lambda \in (0, 1)$.

Proof of Properties

Proof



Actually quite simple

Given

$$y = \lambda x_A + (1 - \lambda) x_B$$

with $\lambda \in (0, 1)$.

Proof of Properties

We know that

$$\begin{aligned}g_k(\mathbf{y}) &= \mathbf{w}^T (\lambda \mathbf{x}_A + (1 - \lambda) \mathbf{x}_B) + w_0 \\&= \lambda \mathbf{w}^T \mathbf{x}_A + \lambda w_0 + (1 - \lambda) \mathbf{w}^T \mathbf{x}_B + (1 - \lambda) w_0 \\&= \lambda g_k(\mathbf{x}_A) + (1 - \lambda) g_k(\mathbf{x}_B) \\&> \lambda g_j(\mathbf{x}_A) + (1 - \lambda) g_j(\mathbf{x}_B) \\&> g_j(\lambda \mathbf{x}_A + (1 - \lambda) \mathbf{x}_B) \\&> g_j(\mathbf{y})\end{aligned}$$

For all $j \neq k$

Proof of Properties

We know that

$$\begin{aligned}g_k(\mathbf{y}) &= \mathbf{w}^T (\lambda \mathbf{x}_A + (1 - \lambda) \mathbf{x}_B) + w_0 \\&= \lambda \mathbf{w}^T \mathbf{x}_A + \lambda w_0 + (1 - \lambda) \mathbf{w}^T \mathbf{x}_B + (1 - \lambda) w_0 \\&= \lambda g_k(\mathbf{x}_A) + (1 - \lambda) g_k(\mathbf{x}_B) \\&> \lambda g_j(\mathbf{x}_A) + (1 - \lambda) g_j(\mathbf{x}_B) \\&> g_j(\lambda \mathbf{x}_A + (1 - \lambda) \mathbf{x}_B) \\&> g_j(\mathbf{y})\end{aligned}$$

For all $j \neq k$

Or...

- \mathbf{y} belongs to an area k defined by the rule!!!
- This area is **Convex** and **Singly Connected** because the definition of \mathbf{y} .

Proof of Properties

We know that

$$\begin{aligned}g_k(\mathbf{y}) &= \mathbf{w}^T (\lambda \mathbf{x}_A + (1 - \lambda) \mathbf{x}_B) + w_0 \\&= \lambda \mathbf{w}^T \mathbf{x}_A + \lambda w_0 + (1 - \lambda) \mathbf{w}^T \mathbf{x}_B + (1 - \lambda) w_0 \\&= \lambda g_k(\mathbf{x}_A) + (1 - \lambda) g_k(\mathbf{x}_B) \\&> \lambda g_j(\mathbf{x}_A) + (1 - \lambda) g_j(\mathbf{x}_B) \\&> g_j(\lambda \mathbf{x}_A + (1 - \lambda) \mathbf{x}_B) \\&> g_j(\mathbf{y})\end{aligned}$$

For all $j \neq k$

Or...

- \mathbf{y} belongs to an area k defined by the rule!!!
- This area is Convex and Singly Connected because the definition of \mathbf{y} .

Proof of Properties

We know that

$$\begin{aligned}g_k(\mathbf{y}) &= \mathbf{w}^T (\lambda \mathbf{x}_A + (1 - \lambda) \mathbf{x}_B) + w_0 \\&= \lambda \mathbf{w}^T \mathbf{x}_A + \lambda w_0 + (1 - \lambda) \mathbf{w}^T \mathbf{x}_B + (1 - \lambda) w_0 \\&= \lambda g_k(\mathbf{x}_A) + (1 - \lambda) g_k(\mathbf{x}_B) \\&> \lambda g_j(\mathbf{x}_A) + (1 - \lambda) g_j(\mathbf{x}_B) \\&> g_j(\lambda \mathbf{x}_A + (1 - \lambda) \mathbf{x}_B) \\&> g_j(\mathbf{y})\end{aligned}$$

For all $j \neq k$

Or...

- \mathbf{y} belongs to an area k defined by the rule!!!
- This area is Convex and Singly Connected because the definition of \mathbf{y} .

Proof of Properties

We know that

$$\begin{aligned}g_k(\mathbf{y}) &= \mathbf{w}^T (\lambda \mathbf{x}_A + (1 - \lambda) \mathbf{x}_B) + w_0 \\&= \lambda \mathbf{w}^T \mathbf{x}_A + \lambda w_0 + (1 - \lambda) \mathbf{w}^T \mathbf{x}_B + (1 - \lambda) w_0 \\&= \lambda g_k(\mathbf{x}_A) + (1 - \lambda) g_k(\mathbf{x}_B) \\&> \lambda g_j(\mathbf{x}_A) + (1 - \lambda) g_j(\mathbf{x}_B) \\&> g_j(\lambda \mathbf{x}_A + (1 - \lambda) \mathbf{x}_B) \\&> g_j(\mathbf{y})\end{aligned}$$

For all $j \neq k$

Or...

- \mathbf{y} belongs to an area k defined by the rule!!!
- This area is Convex and Singly Connected because the definition of \mathbf{y} .

Proof of Properties

We know that

$$\begin{aligned}g_k(\mathbf{y}) &= \mathbf{w}^T (\lambda \mathbf{x}_A + (1 - \lambda) \mathbf{x}_B) + w_0 \\&= \lambda \mathbf{w}^T \mathbf{x}_A + \lambda w_0 + (1 - \lambda) \mathbf{w}^T \mathbf{x}_B + (1 - \lambda) w_0 \\&= \lambda g_k(\mathbf{x}_A) + (1 - \lambda) g_k(\mathbf{x}_B) \\&> \lambda g_j(\mathbf{x}_A) + (1 - \lambda) g_j(\mathbf{x}_B) \\&> g_j(\lambda \mathbf{x}_A + (1 - \lambda) \mathbf{x}_B) \\&> g_j(\mathbf{y})\end{aligned}$$

For all $j \neq k$

Or...

- \mathbf{y} belongs to an area k defined by the rule!!!
- This area is Convex and Singly Connected because the definition of \mathbf{y} .

Proof of Properties

We know that

$$\begin{aligned}g_k(\mathbf{y}) &= \mathbf{w}^T (\lambda \mathbf{x}_A + (1 - \lambda) \mathbf{x}_B) + w_0 \\&= \lambda \mathbf{w}^T \mathbf{x}_A + \lambda w_0 + (1 - \lambda) \mathbf{w}^T \mathbf{x}_B + (1 - \lambda) w_0 \\&= \lambda g_k(\mathbf{x}_A) + (1 - \lambda) g_k(\mathbf{x}_B) \\&> \lambda g_j(\mathbf{x}_A) + (1 - \lambda) g_j(\mathbf{x}_B) \\&> g_j(\lambda \mathbf{x}_A + (1 - \lambda) \mathbf{x}_B) \\&> g_j(\mathbf{y})\end{aligned}$$

For all $j \neq k$

Or...

- \mathbf{y} belongs to an area k defined by the rule!!!
- This area is Convex and Singly Connected because the definition of \mathbf{y} .

Proof of Properties

We know that

$$\begin{aligned}g_k(\mathbf{y}) &= \mathbf{w}^T (\lambda \mathbf{x}_A + (1 - \lambda) \mathbf{x}_B) + w_0 \\&= \lambda \mathbf{w}^T \mathbf{x}_A + \lambda w_0 + (1 - \lambda) \mathbf{w}^T \mathbf{x}_B + (1 - \lambda) w_0 \\&= \lambda g_k(\mathbf{x}_A) + (1 - \lambda) g_k(\mathbf{x}_B) \\&> \lambda g_j(\mathbf{x}_A) + (1 - \lambda) g_j(\mathbf{x}_B) \\&> g_j(\lambda \mathbf{x}_A + (1 - \lambda) \mathbf{x}_B) \\&> g_j(\mathbf{y})\end{aligned}$$

For all $j \neq k$

Or...

- \mathbf{y} belongs to an area k defined by the rule!!!
- This area is Convex and Singly Connected because the definition of \mathbf{y} .

Proof of Properties

We know that

$$\begin{aligned}g_k(\mathbf{y}) &= \mathbf{w}^T (\lambda \mathbf{x}_A + (1 - \lambda) \mathbf{x}_B) + w_0 \\&= \lambda \mathbf{w}^T \mathbf{x}_A + \lambda w_0 + (1 - \lambda) \mathbf{w}^T \mathbf{x}_B + (1 - \lambda) w_0 \\&= \lambda g_k(\mathbf{x}_A) + (1 - \lambda) g_k(\mathbf{x}_B) \\&> \lambda g_j(\mathbf{x}_A) + (1 - \lambda) g_j(\mathbf{x}_B) \\&> g_j(\lambda \mathbf{x}_A + (1 - \lambda) \mathbf{x}_B) \\&> g_j(\mathbf{y})\end{aligned}$$

For all $j \neq k$

Or...

- \mathbf{y} belongs to an area k defined by the rule!!!
- This area is **Convex** and **Singly Connected** because the definition of \mathbf{y} .

However!!!

No so nice properties!!!

- **It limits the power of classification for multi-objective function.**



Cinvestav

How do we train this Linear Machine?

We know that each ω_k class is described by

$$g_k(\mathbf{x}) = \mathbf{w}_k^T \mathbf{x} + w_0 \text{ where } k = 1, \dots, c$$

We then design a single machine

$$g(\mathbf{x}) = W^T \mathbf{x} \quad (24)$$



How do we train this Linear Machine?

We know that each ω_k class is described by

$$g_k(\mathbf{x}) = \mathbf{w}_k^T \mathbf{x} + w_0 \text{ where } k = 1, \dots, c$$

We then design a single machine

$$g(\mathbf{x}) = \mathbf{W}^T \mathbf{x} \quad (24)$$



Where

We have the following

$$\mathbf{W}^T = \begin{pmatrix} 1 & w_{11} & w_{12} & \cdots & w_{1d} \\ 1 & w_{21} & w_{22} & \cdots & w_{2d} \\ 1 & w_{31} & w_{32} & \cdots & w_{3d} \\ \vdots & \vdots & \vdots & & \vdots \\ 1 & w_{c1} & w_{c2} & \cdots & w_{cd} \end{pmatrix} \quad (25)$$

What about the labels?

OK, we know how to do with 2 classes, What about many classes?



Where

We have the following

$$\mathbf{W}^T = \begin{pmatrix} 1 & w_{11} & w_{12} & \cdots & w_{1d} \\ 1 & w_{21} & w_{22} & \cdots & w_{2d} \\ 1 & w_{31} & w_{32} & \cdots & w_{3d} \\ \vdots & \vdots & \vdots & & \vdots \\ 1 & w_{c1} & w_{c2} & \cdots & w_{cd} \end{pmatrix} \quad (25)$$

What about the labels?

OK, we know how to do with 2 classes, What about many classes?



How do we train this Linear Machine?

Use a vector t_i with dimensionality c to identify each element at each class

We have then the following dataset

$$\{\mathbf{x}_i, \mathbf{t}_i\} \text{ for } i = 1, 2, \dots, N$$

We build the following Matrix of Vectors

$$T = \begin{pmatrix} t_1^T \\ t_2^T \\ \vdots \\ t_{N-1}^T \\ t_N^T \end{pmatrix} \quad (26)$$

How do we train this Linear Machine?

Use a vector \mathbf{t}_i with dimensionality c to identify each element at each class

We have then the following dataset

$$\{\mathbf{x}_i, \mathbf{t}_i\} \text{ for } i = 1, 2, \dots, N$$

We build the following Matrix of Vectors

$$\mathbf{T} = \begin{pmatrix} \mathbf{t}_1^T \\ \mathbf{t}_2^T \\ \vdots \\ \mathbf{t}_{N-1}^T \\ \mathbf{t}_N^T \end{pmatrix} \quad (26)$$

Examples for the t_i

Vectors like (One Shot Representation)

$$x_i \neq 0, i \text{ Class} \rightarrow \begin{pmatrix} 0 \\ 0 \\ \vdots \\ 0 \\ 1 \\ 0 \\ \vdots \\ 0 \end{pmatrix}$$

Another possible vector

$$x_i \neq -1, i \text{ Class} \rightarrow \begin{pmatrix} -1 \\ -1 \\ \vdots \\ -1 \\ -1 \\ -1 \\ \vdots \\ -1 \end{pmatrix}$$

Examples for the t_i

Vectors like (One Shot Representation)

$$x_i \neq 0, i \text{ Class} \rightarrow \begin{pmatrix} 0 \\ 0 \\ \vdots \\ 0 \\ 1 \\ 0 \\ \vdots \\ 0 \end{pmatrix}$$

Another possible vector

$$x_i \neq -1, i \text{ Class} \rightarrow \begin{pmatrix} -1 \\ -1 \\ \vdots \\ -1 \\ 1 \\ -1 \\ \vdots \\ -1 \end{pmatrix}$$

Thus, we create the following Matrix

A Matrix containing all the required information

$$XW - T \quad (27)$$

Thus, we create the following Matrix

A Matrix containing all the required information

$$XW - T \quad (27)$$

Where we have the following vector

$$\left[\mathbf{x}_i^T \mathbf{w}_1, \mathbf{x}_i^T \mathbf{w}_2, \mathbf{x}_i^T \mathbf{w}_3, \dots, \mathbf{x}_i^T \mathbf{w}_c \right] \quad (28)$$

Remark: It is the vector result of multiplication of row i of X against W on XW .

Thus, we create the following Matrix

A Matrix containing all the required information

$$XW - T \quad (27)$$

Where we have the following vector

$$\left[\mathbf{x}_i^T \mathbf{w}_1, \mathbf{x}_i^T \mathbf{w}_2, \mathbf{x}_i^T \mathbf{w}_3, \dots, \mathbf{x}_i^T \mathbf{w}_c \right] \quad (28)$$

Remark: It is the vector result of multiplication of row i of X against W on XW .

That is compared to the vector T on T by using the subtraction of vectors

$$e_i = \left[\mathbf{x}_i^T \mathbf{w}_1, \mathbf{x}_i^T \mathbf{w}_2, \mathbf{x}_i^T \mathbf{w}_3, \dots, \mathbf{x}_i^T \mathbf{w}_c \right] - t_i^T \quad (29)$$

Thus, we create the following Matrix

A Matrix containing all the required information

$$\mathbf{XW} - \mathbf{T} \quad (27)$$

Where we have the following vector

$$\left[\mathbf{x}_i^T \mathbf{w}_1, \mathbf{x}_i^T \mathbf{w}_2, \mathbf{x}_i^T \mathbf{w}_3, \dots, \mathbf{x}_i^T \mathbf{w}_c \right] \quad (28)$$

Remark: It is the vector result of multiplication of row i of \mathbf{X} against \mathbf{W} on \mathbf{XW} .

That is compared to the vector \mathbf{t}_i^T on \mathbf{T} by using the subtraction of vectors

$$\mathbf{e}_i = \left[\mathbf{x}_i^T \mathbf{w}_1, \mathbf{x}_i^T \mathbf{w}_2, \mathbf{x}_i^T \mathbf{w}_3, \dots, \mathbf{x}_i^T \mathbf{w}_c \right] - \mathbf{t}_i^T \quad (29)$$

What do we want?

We want the quadratic error

$$\frac{1}{2}e_i^2$$

This specific quadratic errors are at the diagonal of the matrix

$$(XW - T)^T (XW - T)$$

We can use the trace function to generate the desired total error of

$$J(\cdot) = \frac{1}{2} \sum_{i=1}^N e_i^2 \quad (30)$$



What do we want?

We want the quadratic error

$$\frac{1}{2}e_i^2$$

This specific quadratic errors are at the diagonal of the matrix

$$(\mathbf{XW} - \mathbf{T})^T (\mathbf{XW} - \mathbf{T})$$

We can use the trace function to generate the desired total error of

$$J(\cdot) = \frac{1}{2} \sum_{i=1}^N e_i^2 \quad (30)$$



What do we want?

We want the quadratic error

$$\frac{1}{2}e_i^2$$

This specific quadratic errors are at the diagonal of the matrix

$$(\mathbf{XW} - \mathbf{T})^T (\mathbf{XW} - \mathbf{T})$$

We can use the trace function to generate the desired total error of

$$J(\cdot) = \frac{1}{2} \sum_{i=1}^N e_i^2 \quad (30)$$



Then

The trace allows to express the total error

$$J(\mathbf{W}) = \frac{1}{2} \text{Trace} \left\{ (\mathbf{X}\mathbf{W} - \mathbf{T})^T (\mathbf{X}\mathbf{W} - \mathbf{T}) \right\} \quad (31)$$

Thus, we have by the same derivative method

$$\mathbf{W} = (\mathbf{X}^T \mathbf{X})^{-1} \mathbf{X}^T \mathbf{T} = \mathbf{X}^+ \mathbf{T} \quad (32)$$



Then

The trace allows to express the total error

$$J(\mathbf{W}) = \frac{1}{2} \text{Trace} \left\{ (\mathbf{X}\mathbf{W} - \mathbf{T})^T (\mathbf{X}\mathbf{W} - \mathbf{T}) \right\} \quad (31)$$

Thus, we have by the same derivative method

$$\mathbf{W} = (\mathbf{X}^T \mathbf{X})^{-1} \mathbf{X}^T \mathbf{T} = \mathbf{X}^+ \mathbf{T} \quad (32)$$



How do we obtain the discriminant?

Thus, we obtain the discriminant

$$g(\mathbf{x}) = \mathbf{W}^T \mathbf{x} = \mathbf{T}^T (\mathbf{X}^+)^T \mathbf{x} \quad (33)$$



Outline

1

Introduction

- Introduction
- Regression as approximation
- The Simplest Functions
- Splitting the Space
- Defining the Decision Surface
- Properties of the Hyperplane $w^T x + w_0$
- Augmenting the Vector

2

Developing a Solution

- Least Squared Error Procedure
 - The Geometry of a Two-Category Linearly-Separable Case
- The Error Idea
- The Final Error Equation
- Basic Solution
- Multidimensional Solution
- Remember in matrices of 3×3
- What Lives Where?
- Geometric Interpretation
- Solving the Labeling Issue
- Multi-Class Solution
- **Issues with Least Squares!!!**
 - Singularity Notes
 - Problem with Outliers
 - Problem with High Number of Dimensions
- What can be done?
 - Using Statistics to find Important Features
 - What about Numerical Stability?
 - Ridge Regression
- Observation About Eigenvalues

3

Exercises

- Some Stuff for the Lab



Outline

1

Introduction

- Introduction
- Regression as approximation
- The Simplest Functions
- Splitting the Space
- Defining the Decision Surface
- Properties of the Hyperplane $w^T x + w_0$
- Augmenting the Vector

2

Developing a Solution

- Least Squared Error Procedure
 - The Geometry of a Two-Category Linearly-Separable Case
- The Error Idea
- The Final Error Equation
- Basic Solution
- Multidimensional Solution
- Remember in matrices of 3×3
- What Lives Where?
- Geometric Interpretation
- Solving the Labeling Issue
- Multi-Class Solution
- **Issues with Least Squares!!!**
 - **Singularity Notes**
 - Problem with Outliers
 - Problem with High Number of Dimensions
 - What can be done?
 - Using Statistics to find Important Features
 - What about Numerical Stability?
 - Ridge Regression
 - Observation About Eigenvalues

3

Exercises

- Some Stuff for the Lab



Let me show you the covariance matrix

We have in matrix notation

$$S = \frac{1}{N-1} (X - \mathbf{1}\bar{x}^T)^T (X - \mathbf{1}\bar{x}^T)$$

This $X^T X$

It looks a lot like a covariance matrix

Actually, the dependency observed in matrix $X^T X$ between its columns!

- It is the same dependency observed between the features in the data X after the features have been centered by \bar{x} .



Let me show you the covariance matrix

We have in matrix notation

$$S = \frac{1}{N-1} (X - \mathbf{1}\bar{x}^T)^T (X - \mathbf{1}\bar{x}^T)$$

Thus, $X^T X$

It looks a lot like a covariance matrix

Actually, the dependency observed in matrix $X^T X$ between its columns!

- It is the same dependency observed between the features in the data X after the features have been centered by \bar{x} .



Let me show you the covariance matrix

We have in matrix notation

$$S = \frac{1}{N-1} (X - \mathbf{1}\bar{x}^T)^T (X - \mathbf{1}\bar{x}^T)$$

Thus, $X^T X$

It looks a lot like a covariance matrix

Actually, the dependency observed in matrix $X^T X$ between its columns!!!

- It is the same dependency observed between the features in the data X after the features have been centered by \bar{x} .



Thus

We can apply a similar analysis...

- To obtain some of the possible cases that make $X^T X$ singular

A Classical One

- If there is a interdependence between features
 - ▶ Meaning some feature is an exact linear combination of the other features.
 - ▶ The $X^T X$ matrix of the features will be singular.



Cinvestav

Thus

We can apply a similar analysis...

- To obtain some of the possible cases that make $X^T X$ singular

A Classical One

- If there is a interdependence between features
 - ▶ Meaning some feature is an exact linear combination of the other features.
 - ▶ The $X^T X$ matrix of the features will be singular.



Cinvestav

When does this happen?

First

Number of features is equal or greater than the number of samples.

Second

Two or more features sum up to a constant

- For example, $x_2 - 5x_{10} = 0$

Third

Two features are identical or differ merely in mean or variance.



When does this happen?

First

Number of features is equal or greater than the number of samples.

Second

Two or more features sum up to a constant

- For example, $x_2 - 5x_{10} = 0$

Third

Two features are identical or differ merely in mean or variance.



When does this happen?

First

Number of features is equal or greater than the number of samples.

Second

Two or more features sum up to a constant

- For example, $x_2 - 5x_{10} = 0$

Third

Two features are identical or differ merely in mean or variance.



Nevertheless

The least squares coefficients $\hat{\mathbf{w}}$ are not uniquely defined.

- The fitted values $\hat{\mathbf{y}} = \mathbf{X}\hat{\mathbf{w}}$ are still the projection of \mathbf{y} onto the column space of \mathbf{X} .



Additionally

Duplicate observations in a data set

- It will lead the matrix toward singularity.

Caution: scale

- When doing some sort of imputation (Adding missing features), it is always beneficial (from both statistical and mathematical view) to add some noise to the imputed data.

This can happen in the preprocessing phase too

- Be careful.



Additionally

Duplicate observations in a data set

- It will lead the matrix toward singularity.

Cautionary Tale

- When doing some sort of imputation (Adding missing features), it is always beneficial (from both statistical and mathematical view) to add some noise to the imputed data.

This can happen in the preprocessing phase

- Be careful.



Additionally

Duplicate observations in a data set

- It will lead the matrix toward singularity.

Cautionary Tale

- When doing some sort of imputation (Adding missing features), it is always beneficial (from both statistical and mathematical view) to add some noise to the imputed data.

This can happen in the preprocessing phase

- Be careful.



Also

It can happen also that

- $\mathbf{X}^T \mathbf{X}$ could be almost not invertible, making Least Squares numerically unstable.

Statistical consequence

- High variance of predictions.



Also

It can happen also that

- $\mathbf{X}^T \mathbf{X}$ could be almost not invertible, making Least Squares numerically unstable.

Statistical consequence

- High variance of predictions.



When can this happen?

The non-full-rank case occurs

- Most often when one or more qualitative (Categorical Variables/Dummy Variables) inputs are coded in a redundant fashion.

How do we solve this?

- Re-encode or dropping redundant columns in X .

Most regression software packages

- They detect these redundancies and automatically implement some strategies for removing them.



When can this happen?

The non-full-rank case occurs

- Most often when one or more qualitative (Categorical Variables/Dummy Variables) inputs are coded in a redundant fashion.

How do we solve this?

- Re-encode or dropping redundant columns in X .

Model regression software packages

- They detect these redundancies and automatically implement some strategies for removing them.



When can this happen?

The non-full-rank case occurs

- Most often when one or more qualitative (Categorical Variables/Dummy Variables) inputs are coded in a redundant fashion.

How do we solve this?

- Re-encode or dropping redundant columns in X .

Most regression software packages

- They detect these redundancies and automatically implement some strategies for removing them.



Outline

1

Introduction

- Introduction
- Regression as approximation
- The Simplest Functions
- Splitting the Space
- Defining the Decision Surface
- Properties of the Hyperplane $w^T x + w_0$
- Augmenting the Vector

2

Developing a Solution

- Least Squared Error Procedure
 - The Geometry of a Two-Category Linearly-Separable Case
- The Error Idea
- The Final Error Equation
- Basic Solution
- Multidimensional Solution
- Remember in matrices of 3×3
- What Lives Where?
- Geometric Interpretation
- Solving the Labeling Issue
- Multi-Class Solution
- **Issues with Least Squares!!!**
 - Singularity Notes
 - **Problem with Outliers**
 - Problem with High Number of Dimensions
- What can be done?
 - Using Statistics to find Important Features
 - What about Numerical Stability?
 - Ridge Regression
- Observation About Eigenvalues

3

Exercises

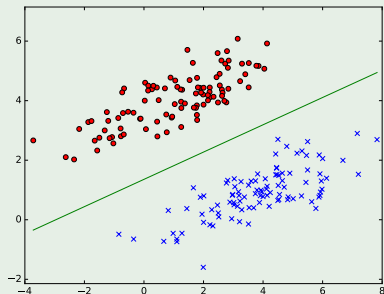
- Some Stuff for the Lab



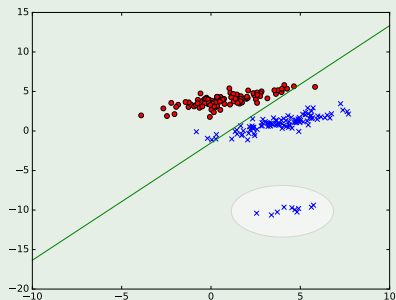
Issues with Least Squares

Problem with Outliers

No Outliers



Outliers



Outline

1

Introduction

- Introduction
- Regression as approximation
- The Simplest Functions
- Splitting the Space
- Defining the Decision Surface
- Properties of the Hyperplane $w^T x + w_0$
- Augmenting the Vector

2

Developing a Solution

- Least Squared Error Procedure
 - The Geometry of a Two-Category Linearly-Separable Case
- The Error Idea
- The Final Error Equation
- Basic Solution
- Multidimensional Solution
- Remember in matrices of 3×3
- What Lives Where?
- Geometric Interpretation
- Solving the Labeling Issue
- Multi-Class Solution
- **Issues with Least Squares!!!**
 - Singularity Notes
 - Problem with Outliers
 - **Problem with High Number of Dimensions**
- What can be done?
 - Using Statistics to find Important Features
 - What about Numerical Stability?
 - Ridge Regression
- Observation About Eigenvalues

3

Exercises

- Some Stuff for the Lab



Cinvestav

Problems with a High Number of Dimensions

In Many Modern Problems

- Many dimensions/features/predictors (possibly thousands).



Problems with a High Number of Dimensions

In Many Modern Problems

- Many dimensions/features/predictors (possibly thousands).

Only a few of these may be important

- It needs some form of feature selection.
- Possible some type of regularization.



Problems with a High Number of Dimensions

In Many Modern Problems

- Many dimensions/features/predictors (possibly thousands).

Only a few of these may be important

- It needs some form of feature selection.
- Possible some type of regularization.

Why?

- Least Square Error Regression treats all dimensions equally.
- Relevant dimensions might be averaged with irrelevant ones.



Problems with a High Number of Dimensions

In Many Modern Problems

- Many dimensions/features/predictors (possibly thousands).

Only a few of these may be important

- It needs some form of feature selection.
- Possible some type of regularization.

Why?

- Least Square Error Regression treats all dimensions equally.
- Relevant dimensions might be averaged with irrelevant ones.



Problems with a High Number of Dimensions

In Many Modern Problems

- Many dimensions/features/predictors (possibly thousands).

Only a few of these may be important

- It needs some form of feature selection.
- Possible some type of regularization.

Why?

- Least Square Error Regression treats all dimensions equally.
- Relevant dimensions might be averaged with irrelevant ones.



Outline

1

Introduction

- Introduction
- Regression as approximation
- The Simplest Functions
- Splitting the Space
- Defining the Decision Surface
- Properties of the Hyperplane $w^T x + w_0$
- Augmenting the Vector

2

Developing a Solution

- Least Squared Error Procedure
 - The Geometry of a Two-Category Linearly-Separable Case
- The Error Idea
- The Final Error Equation
- Basic Solution
- Multidimensional Solution
- Remember in matrices of 3×3
- What Lives Where?
- Geometric Interpretation
- Solving the Labeling Issue
- Multi-Class Solution
- Issues with Least Squares!!!
 - Singularity Notes
 - Problem with Outliers
 - Problem with High Number of Dimensions
- **What can be done?**
 - Using Statistics to find Important Features
 - What about Numerical Stability?
 - Ridge Regression
- Observation About Eigenvalues

3

Exercises

- Some Stuff for the Lab



Outline

1

Introduction

- Introduction
- Regression as approximation
- The Simplest Functions
- Splitting the Space
- Defining the Decision Surface
- Properties of the Hyperplane $w^T x + w_0$
- Augmenting the Vector

2

Developing a Solution

- Least Squared Error Procedure
 - The Geometry of a Two-Category Linearly-Separable Case
- The Error Idea
- The Final Error Equation
- Basic Solution
- Multidimensional Solution
- Remember in matrices of 3×3
- What Lives Where?
- Geometric Interpretation
- Solving the Labeling Issue
- Multi-Class Solution
- Issues with Least Squares!!!
 - Singularity Notes
 - Problem with Outliers
 - Problem with High Number of Dimensions
- What can be done?
 - **Using Statistics to find Important Features**
 - What about Numerical Stability?
 - Ridge Regression
- Observation About Eigenvalues

3

Exercises

- Some Stuff for the Lab



We will start using some statistics

We want to obtain sampling properties for $\hat{\mathbf{w}}$

For this remember:

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

For this assume:

- The observations y_i are uncorrelated and have constant variance σ^2 .
- The x_i are fixed = not random.



We will start using some statistics

We want to obtain sampling properties for \hat{w}

For this remember:

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

For this assume,

- The observations y_i are uncorrelated and have constant variance σ^2 .
- The x_i are fixed = not random.



Then, we have the variance-covariance matrix

We have

$$\text{Var}(\hat{\mathbf{w}}) = \text{Var} \left[\left(\mathbf{X}^T \mathbf{X} \right)^{-1} \mathbf{X}^T \mathbf{y} \right]$$

We have the following equivalence:

$$\text{Var}(A\mathbf{y}) = A \text{Var}(\mathbf{y}) A^T$$



Then, we have the variance-covariance matrix

We have

$$\text{Var}(\hat{\mathbf{w}}) = \text{Var} \left[\left(\mathbf{X}^T \mathbf{X} \right)^{-1} \mathbf{X}^T \mathbf{y} \right]$$

We have the following equivalence

$$\text{Var}(\mathbf{A}\mathbf{y}) = \mathbf{A} \text{Var}(\mathbf{y}) \mathbf{A}^T$$



Therefore

Something Notable

$$\begin{aligned} \text{Var} \left[\left(\mathbf{X}^T \mathbf{X} \right)^{-1} \mathbf{X}^T \mathbf{y} \right] &= \left(\mathbf{X}^T \mathbf{X} \right)^{-1} \mathbf{X}^T \text{Var}(\mathbf{y}) \mathbf{X} \left(\mathbf{X}^T \mathbf{X} \right)^{-1} \\ &= \left(\mathbf{X}^T \mathbf{X} \right)^{-1} \mathbf{X}^T \sigma^2 \mathbf{I} \mathbf{X} \left(\mathbf{X}^T \mathbf{X} \right)^{-1} \\ &= \sigma^2 \left(\mathbf{X}^T \mathbf{X} \right)^{-1} \end{aligned}$$

Therefore

Something Notable

$$\begin{aligned} \text{Var} \left[\left(\mathbf{X}^T \mathbf{X} \right)^{-1} \mathbf{X}^T \mathbf{y} \right] &= \left(\mathbf{X}^T \mathbf{X} \right)^{-1} \mathbf{X}^T \text{Var}(\mathbf{y}) \mathbf{X} \left(\mathbf{X}^T \mathbf{X} \right)^{-1} \\ &= \left(\mathbf{X}^T \mathbf{X} \right)^{-1} \mathbf{X}^T \sigma^2 \mathbf{I} \mathbf{X} \left(\mathbf{X}^T \mathbf{X} \right)^{-1} \\ &= \sigma^2 \left(\mathbf{X}^T \mathbf{X} \right)^{-1} \end{aligned}$$

Given that

$$\text{Var}(\mathbf{y}) = \begin{bmatrix} \text{Var}(y_1) & \text{Cov}(y_1, y_2) & \cdots & \text{Cov}(y_1, y_N) \\ \text{Cov}(y_2, y_1) & \cdots & \text{Var}(y_2) & \cdots & \text{Cov}(y_2, y_N) \\ \vdots & \vdots & \vdots & \ddots & \vdots \\ \text{Cov}(y_N, y_1) & \text{Cov}(y_N, y_2) & \cdots & \text{Var}(y_N) \end{bmatrix} = \begin{bmatrix} \sigma^2 & 0 & \cdots & 0 \\ 0 & \sigma^2 & \cdots & 0 \\ \vdots & \vdots & \ddots & \vdots \\ 0 & 0 & \cdots & \sigma^2 \end{bmatrix}$$

Therefore

Something Notable

$$\begin{aligned}\text{Var} \left[\left(\mathbf{X}^T \mathbf{X} \right)^{-1} \mathbf{X}^T \mathbf{y} \right] &= \left(\mathbf{X}^T \mathbf{X} \right)^{-1} \mathbf{X}^T \text{Var}(\mathbf{y}) \mathbf{X} \left(\mathbf{X}^T \mathbf{X} \right)^{-1} \\ &= \left(\mathbf{X}^T \mathbf{X} \right)^{-1} \mathbf{X}^T \sigma^2 \mathbf{I} \mathbf{X} \left(\mathbf{X}^T \mathbf{X} \right)^{-1} \\ &= \sigma^2 \left(\mathbf{X}^T \mathbf{X} \right)^{-1}\end{aligned}$$

Given that

$$\text{Var}(\mathbf{y}) = \begin{bmatrix} \text{Var}(y_1) & \text{Cov}(y_1, y_2) & \cdots & \text{Cov}(y_1, y_N) \\ \text{Cov}(y_2, y_1) & \cdots & \text{Var}(y_2) & \cdots & \text{Cov}(y_2, y_N) \\ \vdots & \vdots & \vdots & \ddots & \vdots \\ \text{Cov}(y_N, y_1) & \text{Cov}(y_N, y_2) & \cdots & \text{Var}(y_N) \end{bmatrix} = \begin{bmatrix} \sigma^2 & 0 & \cdots & 0 \\ 0 & \sigma^2 & \cdots & 0 \\ \vdots & \vdots & \ddots & \vdots \\ 0 & 0 & \cdots & \sigma^2 \end{bmatrix}$$

Therefore

Something Notable

$$\begin{aligned} \text{Var} \left[\left(\mathbf{X}^T \mathbf{X} \right)^{-1} \mathbf{X}^T \mathbf{y} \right] &= \left(\mathbf{X}^T \mathbf{X} \right)^{-1} \mathbf{X}^T \text{Var}(\mathbf{y}) \mathbf{X} \left(\mathbf{X}^T \mathbf{X} \right)^{-1} \\ &= \left(\mathbf{X}^T \mathbf{X} \right)^{-1} \mathbf{X}^T \sigma^2 \mathbf{I} \mathbf{X} \left(\mathbf{X}^T \mathbf{X} \right)^{-1} \\ &= \sigma^2 \left(\mathbf{X}^T \mathbf{X} \right)^{-1} \end{aligned}$$

Given that

$$\text{Var}(\mathbf{y}) = \begin{bmatrix} \text{Var}(y_1) & \text{Cov}(y_1, y_2) & \cdots & \text{Cov}(y_1, y_N) \\ \text{Cov}(y_2, y_1) & \cdots \text{Var}(y_2) & \cdots & \text{Cov}(y_2, y_N) \\ \vdots & \vdots & \ddots & \vdots \\ \text{Cov}(y_N, y_1) & \text{Cov}(y_N, y_2) & \cdots & \text{Var}(y_N) \end{bmatrix} = \begin{bmatrix} \sigma^2 & 0 & \cdots & 0 \\ 0 & \sigma^2 & \cdots & 0 \\ \vdots & \vdots & \ddots & \vdots \\ 0 & 0 & \cdots & \sigma^2 \end{bmatrix}$$

Thus

Typically, we can use the following unbiased estimator

$$\hat{\sigma}^2 = \frac{1}{N - d - 1} \sum_{i=1}^N (y_i - \hat{y}_i)$$

- Which is an unbiased estimator $E[\hat{\sigma}^2] = \sigma^2$.

If we have the following relation

$$Y = E(Y|X_1, X_2, \dots, X_d) + \epsilon$$

where

- $\epsilon \sim N(0, \sigma^2)$



Thus

Typically, we can use the following unbiased estimator

$$\hat{\sigma}^2 = \frac{1}{N - d - 1} \sum_{i=1}^N (y_i - \hat{y}_i)$$

- Which is an unbiased estimator $E[\hat{\sigma}^2] = \sigma^2$.

If we have the following relation

$$Y = E(Y|X_1, X_2, \dots, X_d) + \epsilon$$

where

- $\epsilon \sim N(0, \sigma^2)$



Thus

Typically, we can use the following unbiased estimator

$$\hat{\sigma}^2 = \frac{1}{N - d - 1} \sum_{i=1}^N (y_i - \hat{y}_i)$$

- Which is an unbiased estimator $E [\hat{\sigma}^2] = \sigma^2$.

If we have the following relation

$$Y = E (Y|X_1, X_2, \dots, X_d) + \epsilon$$

Where

- $\epsilon \sim N (0, \sigma^2)$



Then

We have

$$\hat{\mathbf{w}} \sim N\left(\mathbf{w}, \sigma^2 (\mathbf{X}^T \mathbf{X})^{-1}\right)$$

Thus, we can be a little bit smart:

$$H_0 : w_j = 0$$

$$H_1 : w_j \neq 0$$

To test for Hypothesis $w_j = 0$, we get the following t -score:

$$z_j = \frac{\hat{w}_j - w_j}{\hat{\sigma} \sqrt{v_j}} = \frac{\hat{w}_j}{\hat{\sigma} \sqrt{v_j}} \text{ with } v_j \text{ the } j^{\text{th}} \text{ diagonal element at } (\mathbf{X}^T \mathbf{X})^{-1}$$

Then

We have

$$\hat{\mathbf{w}} \sim N\left(\mathbf{w}, \sigma^2 (\mathbf{X}^T \mathbf{X})^{-1}\right)$$

Thus, we can be a little bit smart

$$H_0 : w_j = 0$$

$$H_1 : w_j \neq 0$$

To test for Hypothesis $w_j = 0$, we get the following t -score

$$z_j = \frac{\hat{w}_j - w_j}{\hat{\sigma} \sqrt{v_j}} = \frac{\hat{w}_j}{\hat{\sigma} \sqrt{v_j}} \text{ with } v_j \text{ the } j^{\text{th}} \text{ diagonal element at } (\mathbf{X}^T \mathbf{X})^{-1}$$

Then

We have

$$\hat{\mathbf{w}} \sim N\left(\mathbf{w}, \sigma^2 (\mathbf{X}^T \mathbf{X})^{-1}\right)$$

Thus, we can be a little bit smart

$$H_0 : w_j = 0$$

$$H_1 : w_j \neq 0$$

To test for Hypothesis $w_j = 0$, we get the following z -score

$$z_j = \frac{\hat{w}_j - w_j}{\hat{\sigma} \sqrt{v_j}} = \frac{\hat{w}_j}{\hat{\sigma} \sqrt{v_j}} \text{ with } v_j \text{ the } j^{\text{th}} \text{ diagonal element at } (\mathbf{X}^T \mathbf{X})^{-1}$$

Therefore

$z_j \sim t_{N-d-1}$ a t-student distribution

- Therefore, a large(absolute) value of z_j will lead to rejection of the Null Hypothesis

Therefore

$z_j \sim t_{N-d-1}$ a t-student distribution

- Therefore, a large(absolute) value of z_j will lead to rejection of the Null Hypothesis

Therefore

You can use the simple rule:

- Accept H_0 remove the feature
- Reject H_0 keep the feature

Therefore

$z_j \sim t_{N-d-1}$ a t-student distribution

- Therefore, a large(absolute) value of z_j will lead to rejection of the Null Hypothesis

Therefore

You can use the simple rule:

- Accept H_0 remove the feature
- Reject H_0 keep the feature

However

There are still more techniques for feature selection quite more advanced...

Therefore

$z_j \sim t_{N-d-1}$ a t-student distribution

- Therefore, a large(absolute) value of z_j will lead to rejection of the Null Hypothesis

Therefore

You can use the simple rule:

- Accept H_0 remove the feature
- Reject H_0 keep the feature

However

There are still more techniques for feature selection quite more advanced...

Therefore

$z_j \sim t_{N-d-1}$ a t-student distribution

- Therefore, a large(absolute) value of z_j will lead to rejection of the Null Hypothesis

Therefore

You can use the simple rule:

- Accept H_0 remove the feature
- Reject H_0 keep the feature

However

There are still more techniques for feature selection quite more advanced...

Outline

1

Introduction

- Introduction
- Regression as approximation
- The Simplest Functions
- Splitting the Space
- Defining the Decision Surface
- Properties of the Hyperplane $w^T x + w_0$
- Augmenting the Vector

2

Developing a Solution

- Least Squared Error Procedure
 - The Geometry of a Two-Category Linearly-Separable Case
- The Error Idea
- The Final Error Equation
- Basic Solution
- Multidimensional Solution
- Remember in matrices of 3×3
- What Lives Where?
- Geometric Interpretation
- Solving the Labeling Issue
- Multi-Class Solution
- Issues with Least Squares!!!
 - Singularity Notes
 - Problem with Outliers
 - Problem with High Number of Dimensions
- **What can be done?**
 - Using Statistics to find Important Features
 - **What about Numerical Stability?**
 - Ridge Regression
- Observation About Eigenvalues

3

Exercises

- Some Stuff for the Lab



What to Do About Numerical Stability?

Definition

- A matrix which is not invertible is also called a **singular** matrix.
- A matrix which is invertible (not singular) is called **regular**.

What to Do About Numerical Stability?

Definition

- A matrix which is not invertible is also called a **singular** matrix.
- A matrix which is invertible (not singular) is called **regular**.

What is the Meaning?

Imagine the following in \mathbb{R}^3

$$A = \begin{pmatrix} a_{11} & a_{12} & a_{13} \\ a_{21} & a_{22} & a_{23} \\ a_{31} & a_{32} & a_{33} \end{pmatrix}$$

What to Do About Numerical Stability?

Definition

- A matrix which is not invertible is also called a **singular** matrix.
- A matrix which is invertible (not singular) is called **regular**.

What is the Meaning?

Imagine the following in \mathbb{R}^3

$$A = \begin{pmatrix} a_{11} & a_{12} & a_{13} \\ a_{21} & a_{22} & a_{23} \\ a_{31} & a_{32} & a_{33} \end{pmatrix}$$

Given that the columns are vectors

They span a subspace for those column vectors in \mathbb{R}^3

$$\text{span} \left\{ \begin{pmatrix} a_{11} \\ a_{21} \\ a_{31} \end{pmatrix}, \begin{pmatrix} a_{12} \\ a_{22} \\ a_{32} \end{pmatrix}, \begin{pmatrix} a_{13} \\ a_{23} \\ a_{33} \end{pmatrix} \right\}$$

What to Do About Numerical Stability?

Definition

- A matrix which is not invertible is also called a **singular** matrix.
- A matrix which is invertible (not singular) is called **regular**.

What is the Meaning?

Imagine the following in \mathbb{R}^3

$$A = \begin{pmatrix} a_{11} & a_{12} & a_{13} \\ a_{21} & a_{22} & a_{23} \\ a_{31} & a_{32} & a_{33} \end{pmatrix}$$

Given that the columns are vectors

They span a subspace for those column vectors in \mathbb{R}^3

$$\text{span} \left\{ \begin{pmatrix} a_{11} \\ a_{21} \\ a_{31} \end{pmatrix}, \begin{pmatrix} a_{12} \\ a_{22} \\ a_{32} \end{pmatrix}, \begin{pmatrix} a_{13} \\ a_{23} \\ a_{33} \end{pmatrix} \right\}$$

Relation with the Rank

If a matrix is singular

Its Rank is less than 3, i.e :

- The subspace is squashed into a plane.
- The subspace is squashed into a line.
- The subspace in the WORST CASE into a point.



Relation with the Rank

If a matrix is singular

Its Rank is less than 3, i.e :

- 1 The subspace is squashed into a plane.
- 2 The subspace is squashed into a line.
- 3 The subspace in the WORST CASE into a point.



Relation with the Rank

If a matrix is singular

Its Rank is less than 3, i.e :

- 1 The subspace is squashed into a plane.
- 2 The subspace is squashed into a line.
- 3 The subspace in the WORST CASE into a point.



Relation with the Rank

If a matrix is singular

Its Rank is less than 3, i.e :

- 1 The subspace is squashed into a plane.
- 2 The subspace is squashed into a line.
- 3 The subspace in the WORST CASE into a point.



Remember

That, we have

$$\mathbf{v} = \lambda_1 \begin{pmatrix} a_{11} \\ a_{21} \\ a_{31} \end{pmatrix} + \lambda_2 \begin{pmatrix} a_{12} \\ a_{22} \\ a_{32} \end{pmatrix} + \lambda_3 \begin{pmatrix} a_{13} \\ a_{23} \\ a_{33} \end{pmatrix}$$

Thus, if for example, the matrix projects into a plane

$$\begin{aligned} \mathbf{v} &= \lambda_1 \begin{pmatrix} a_{11} \\ a_{21} \\ a_{31} \end{pmatrix} + \lambda_2 \left[\alpha_1 \begin{pmatrix} a_{11} \\ a_{21} \\ a_{31} \end{pmatrix} + \alpha_2 \begin{pmatrix} a_{13} \\ a_{23} \\ a_{33} \end{pmatrix} \right] + \lambda_3 \begin{pmatrix} a_{13} \\ a_{23} \\ a_{33} \end{pmatrix} \\ &= c_1 \begin{pmatrix} a_{11} \\ a_{21} \\ a_{31} \end{pmatrix} + c_2 \begin{pmatrix} a_{13} \\ a_{23} \\ a_{33} \end{pmatrix} \quad \text{with } c_1 = \lambda_1 + \alpha_1 \lambda_2, c_2 = \alpha_2 \lambda_2 + \lambda_3 \end{aligned}$$

Remember

That, we have

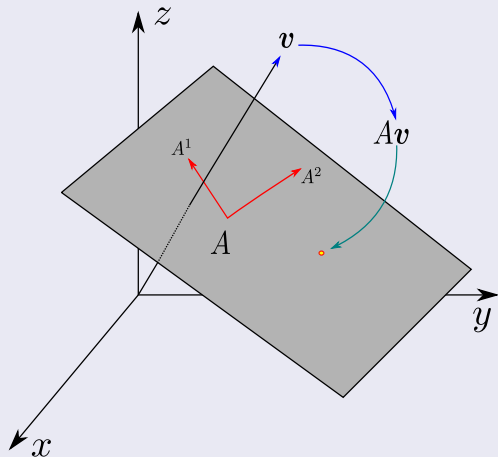
$$\mathbf{v} = \lambda_1 \begin{pmatrix} a_{11} \\ a_{21} \\ a_{31} \end{pmatrix} + \lambda_2 \begin{pmatrix} a_{12} \\ a_{22} \\ a_{32} \end{pmatrix} + \lambda_3 \begin{pmatrix} a_{13} \\ a_{23} \\ a_{33} \end{pmatrix}$$

Thus, if for example, the matrix projects into a plane

$$\begin{aligned} \mathbf{v} &= \lambda_1 \begin{pmatrix} a_{11} \\ a_{21} \\ a_{31} \end{pmatrix} + \lambda_2 \left[\alpha_1 \begin{pmatrix} a_{11} \\ a_{21} \\ a_{31} \end{pmatrix} + \alpha_2 \begin{pmatrix} a_{13} \\ a_{23} \\ a_{33} \end{pmatrix} \right] + \lambda_3 \begin{pmatrix} a_{13} \\ a_{23} \\ a_{33} \end{pmatrix} \\ &= c_1 \begin{pmatrix} a_{11} \\ a_{21} \\ a_{31} \end{pmatrix} + c_2 \begin{pmatrix} a_{13} \\ a_{23} \\ a_{33} \end{pmatrix} \quad \text{with } c_1 = \lambda_1 + \alpha_1 \lambda_2, c_2 = \alpha_2 \lambda_2 + \lambda_3 \end{aligned}$$

For Example

We have a squashing into a plane



Computational Intuition

First Intuition

A singular matrix maps an entire linear subspace into a single point.

Second Intuition

If a matrix maps points far away from each other to points very close to each other, it almost behaves like a singular matrix.



Computational Intuition

First Intuition

A singular matrix maps an entire linear subspace into a single point.

Second Intuitions

If a matrix maps points far away from each other to points very close to each other, it almost behaves like a singular matrix.



Thus

Mapping is related to the eigenvalues!!!

- **Large positive eigenvalues \Rightarrow the mapping is large!!!**

• Small positive eigenvalues \Rightarrow the mapping is small!!!



Cinvestav

Thus

Mapping is related to the eigenvalues!!!

- **Large positive eigenvalues** \Rightarrow **the mapping is large!!!**
- **Small positive eigenvalues** \Rightarrow **the mapping is small!!!**



Cinvestav

There is a statement to support this

All this comes from the following statement

A positive semi-definite matrix A is singular \iff smallest eigenvalue is 0

Consequences for Statistics

If a statistical prediction involves the inverse of an almost-singular matrix, the predictions become unreliable (high variance).



Cinvestav

There is a statement to support this

All this comes from the following statement

A positive semi-definite matrix A is singular \iff smallest eigenvalue is 0

Consequence for Statistics

If a statistical prediction involves the inverse of an almost-singular matrix, the predictions become unreliable (high variance).



Outline

1

Introduction

- Introduction
- Regression as approximation
- The Simplest Functions
- Splitting the Space
- Defining the Decision Surface
- Properties of the Hyperplane $w^T x + w_0$
- Augmenting the Vector

2

Developing a Solution

- Least Squared Error Procedure
 - The Geometry of a Two-Category Linearly-Separable Case
- The Error Idea
- The Final Error Equation
- Basic Solution
- Multidimensional Solution
- Remember in matrices of 3×3
- What Lives Where?
- Geometric Interpretation
- Solving the Labeling Issue
- Multi-Class Solution
- Issues with Least Squares!!!
 - Singularity Notes
 - Problem with Outliers
 - Problem with High Number of Dimensions
- **What can be done?**
 - Using Statistics to find Important Features
 - What about Numerical Stability?
 - **Ridge Regression**
- Observation About Eigenvalues

3

Exercises

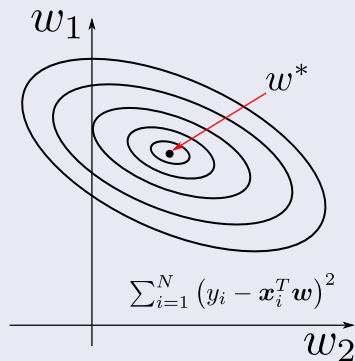
- Some Stuff for the Lab



What can be done?

What could be the problem?

- Imagine that you finish with an over-fitting at the optimal w^*



Overfitting?

Basically (Intuition)

- $x_i^T w^* \approx y_i$

Then

- You are quite good with the training data
- But Really bad with the validation and testing data

We need to pull the optimal in some way!!!

IDEAS?



Overfitting?

Basically (Intuition)

- $x_i^T w^* \approx y_i$

Then

- You are quite good with the training data
- But Really bad with the validation and testing data

We need to pull the optimal in some way!!!

IDEAS?



Overfitting?

Basically (Intuition)

- $x_i^T w^* \approx y_i$

Then

- You are quite good with the training data
- But Really bad with the validation and testing data

We need to pull the optimal in some way!!!

IDEAS?



How do we integrate this solution to the Least Squared Error Solution?

We modify it by adding an extra parameter and tweak the λ

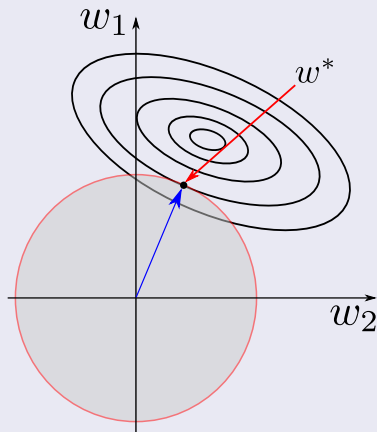
$$\sum_{i=1}^N (y_i - \mathbf{x}_i^T \mathbf{w})^2 + \lambda \sum_{i=1}^{d+1} w_i^2 \quad (34)$$



How do we integrate this solution to the Least Squared Error Solution?

Geometrically Equivalent to pulling away the optimal, it is known as Ridge Regression

$$\sum_{i=1}^N (y_i - \mathbf{x}_i^T \mathbf{w})^2 + \lambda \sum_{i=1}^{d+1} w_i^2$$



Outline

1

Introduction

- Introduction
- Regression as approximation
- The Simplest Functions
- Splitting the Space
- Defining the Decision Surface
- Properties of the Hyperplane $w^T x + w_0$
- Augmenting the Vector

2

Developing a Solution

- Least Squared Error Procedure
 - The Geometry of a Two-Category Linearly-Separable Case
- The Error Idea
- The Final Error Equation
- Basic Solution
- Multidimensional Solution
- Remember in matrices of 3×3
- What Lives Where?
- Geometric Interpretation
- Solving the Labeling Issue
- Multi-Class Solution
- Issues with Least Squares!!!
 - Singularity Notes
 - Problem with Outliers
 - Problem with High Number of Dimensions
- What can be done?
 - Using Statistics to find Important Features
 - What about Numerical Stability?
 - Ridge Regression
- Observation About Eigenvalues

3

Exercises

- Some Stuff for the Lab



Something quite interesting

The w_i in the vector \mathbf{w}^* are related to the eigenvalues in $\mathbf{X}^T \mathbf{X}$

- Thus, we can tweak the eigenvalues to obtain a similar effect than in the Ridge Regression

$$\sum_{i=1}^N (y_i - \mathbf{x}_i^T \mathbf{w})^2 + \lambda \sum_{i=1}^{d+1} w_i^2 \quad (35)$$

It is equivalent to avoid eigenvalues to become zero!!!

Thus, we can do the following given that $\mathbf{X}^T \mathbf{X}$ is positive definite

Assume that $\xi_1, \xi_2, \dots, \xi_{d+1}$ are eigenvectors of $\mathbf{X}^T \mathbf{X}$ with eigenvalues $\lambda_1, \lambda_2, \dots, \lambda_{d+1}$

We have

$$\left(\mathbf{X}^T \mathbf{X}\right) \xi_i = \lambda_i \xi_i \text{ for all } i = 1, \dots, d+1 \quad (36)$$

Given that $\mathbf{X}^T \mathbf{X}$ is singular, some λ_i is equal to 0.

Very Simple, add a convenient λ

$$\left(\mathbf{X}^T \mathbf{X} + \lambda \mathbf{I}\right) \xi_i = (\lambda_i + \lambda) \xi_i \quad (37)$$

i.e. $\lambda_i + \lambda$ is an eigenvalue for $\left(\mathbf{X}^T \mathbf{X} + \lambda \mathbf{I}\right)$.

It is equivalent to avoid eigenvalues to become zero!!!

Thus, we can do the following given that $\mathbf{X}^T \mathbf{X}$ is positive definite

Assume that $\xi_1, \xi_2, \dots, \xi_{d+1}$ are eigenvectors of $\mathbf{X}^T \mathbf{X}$ with eigenvalues $\lambda_1, \lambda_2, \dots, \lambda_{d+1}$

We have

$$\left(\mathbf{X}^T \mathbf{X}\right) \xi_i = \lambda_i \xi_i \text{ for all } i = 1, \dots, d + 1 \quad (36)$$

Given that $\mathbf{X}^T \mathbf{X}$ is singular, some λ_i is equal to 0.

Very Simple and convenient

$$\left(\mathbf{X}^T \mathbf{X} + \lambda I\right) \xi_i = (\lambda_i + \lambda) \xi_i \quad (37)$$

i.e. $\lambda_i + \lambda$ is an eigenvalue for $\left(\mathbf{X}^T \mathbf{X} + \lambda I\right)$.

It is equivalent to avoid eigenvalues to become zero!!!

Thus, we can do the following given that $\mathbf{X}^T \mathbf{X}$ is positive definite

Assume that $\xi_1, \xi_2, \dots, \xi_{d+1}$ are eigenvectors of $\mathbf{X}^T \mathbf{X}$ with eigenvalues $\lambda_1, \lambda_2, \dots, \lambda_{d+1}$

We have

$$\left(\mathbf{X}^T \mathbf{X}\right) \xi_i = \lambda_i \xi_i \text{ for all } i = 1, \dots, d + 1 \quad (36)$$

Given that $\mathbf{X}^T \mathbf{X}$ is singular, some λ_i is equal to 0.

Very Simple, add a convenient λ

$$\left(\mathbf{X}^T \mathbf{X} + \lambda I\right) \xi_i = (\lambda_i + \lambda) \xi_i \quad (37)$$

i.e. $\lambda_i + \lambda$ is an eigenvalue for $\left(\mathbf{X}^T \mathbf{X} + \lambda I\right)$.

What does this mean?

Something Notable

You can control the singularity by detecting the smallest eigenvalue.

Thus

We add an appropriate tuning value λ .



Cinvestav

What does this mean?

Something Notable

You can control the singularity by detecting the smallest eigenvalue.

Thus

We add an appropriate tuning value λ .



Ridge Regression

Ridge Regression

It tries to make least squares more robust if $X^T X$ is almost singular.



Ridge Regression

Ridge Regression

It tries to make least squares more robust if $\mathbf{X}^T \mathbf{X}$ is almost singular.

Process

- 1 Find the eigenvalues of $\mathbf{X}^T \mathbf{X}$
 - 2 If all of them are big enough than zero we are fine!!!
 - 3 Find the smallest one, then tune if necessary.
 - 4 Build $\hat{\mathbf{w}}^{Ridge} = (\mathbf{X}^T \mathbf{X} + \lambda \mathbf{I})^{-1} \mathbf{X}^T \mathbf{y}$.



Ridge Regression

Ridge Regression

It tries to make least squares more robust if $\mathbf{X}^T \mathbf{X}$ is almost singular.

Process

- 1 Find the eigenvalues of $\mathbf{X}^T \mathbf{X}$
- 2 If all of them are bigger enough than zero we are fine!!!

3 Find the smallest one, then tune if necessary.

4 Build $\hat{\mathbf{w}}^{Ridge} = (\mathbf{X}^T \mathbf{X} + \lambda \mathbf{I})^{-1} \mathbf{X}^T \mathbf{y}$.



Ridge Regression

Ridge Regression

It tries to make least squares more robust if $\mathbf{X}^T \mathbf{X}$ is almost singular.

Process

- 1 Find the eigenvalues of $\mathbf{X}^T \mathbf{X}$
- 2 If all of them are bigger enough than zero we are fine!!!
- 3 Find the smallest one, then tune if necessary.

4 Build $\hat{\mathbf{w}}^{Ridge} = (\mathbf{X}^T \mathbf{X} + \lambda \mathbf{I})^{-1} \mathbf{X}^T \mathbf{y}$.



Ridge Regression

Ridge Regression

It tries to make least squares more robust if $\mathbf{X}^T \mathbf{X}$ is almost singular.

Process

- 1 Find the eigenvalues of $\mathbf{X}^T \mathbf{X}$
- 2 If all of them are bigger enough than zero we are fine!!!
- 3 Find the smallest one, then tune if necessary.
- 4 Build $\hat{\mathbf{w}}^{Ridge} = (\mathbf{X}^T \mathbf{X} + \lambda \mathbf{I})^{-1} \mathbf{X}^T \mathbf{y}$.



Outline

1

Introduction

- Introduction
- Regression as approximation
- The Simplest Functions
- Splitting the Space
- Defining the Decision Surface
- Properties of the Hyperplane $w^T x + w_0$
- Augmenting the Vector

2

Developing a Solution

- Least Squared Error Procedure
 - The Geometry of a Two-Category Linearly-Separable Case
- The Error Idea
- The Final Error Equation
- Basic Solution
- Multidimensional Solution
- Remember in matrices of 3×3
- What Lives Where?
- Geometric Interpretation
- Solving the Labeling Issue
- Multi-Class Solution
- Issues with Least Squares!!!
 - Singularity Notes
 - Problem with Outliers
 - Problem with High Number of Dimensions
- What can be done?
 - Using Statistics to find Important Features
 - What about Numerical Stability?
 - Ridge Regression
- Observation About Eigenvalues

3

Exercises

- Some Stuff for the Lab



Exercises

Duda and Hart

Chapter 5

- 1, 3, 4, 7, 13, 17

Bishop

Chapter 4

- 4.1, 4.4, 4.7.

Frank Buss

Chapter 3 - Problems

- Ex 3.5
- Ex 3.6



Exercises

Duda and Hart

Chapter 5

- 1, 3, 4, 7, 13, 17

Bishop

Chapter 4

- 4.1, 4.4, 4.7,

Exercise Problems

Chapter 3 - Problems

- Ex 3.5
- Ex 3.6

Exercises

Duda and Hart

Chapter 5

- 1, 3, 4, 7, 13, 17

Bishop

Chapter 4

- 4.1, 4.4, 4.7,

Hastie-Tibishirani

Chapter 3 - Problems

- Ex 3.5
- Ex 3.6



Theodoridis

Chapter 3 - Problems

