

The Data Science Cycle

Unsupervised Learning - Clustering

Andres Mendez-Vazquez

September 21, 2016

Outline

1 Introduction

- The Simplest Functions
- Splitting the Space
- The Decision Surface
- Minimum Squared Error Procedure
- The Error Idea
- The Final Error Equation
- The Data Matrix
- Issues with Least Squares!!!

Outline

1 Introduction

- The Simplest Functions
 - Splitting the Space
 - The Decision Surface
 - Minimum Squared Error Procedure
 - The Error Idea
 - The Final Error Equation
 - The Data Matrix
 - Issues with Least Squares!!!

What is it?

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

Here, we have an hyperplane as a model:

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

in the case of \mathbb{R}^2 :

We have the following function:

$$g(x) = w_1 x_1 + w_2 x_2 + w_0 \quad (2)$$

What is it?

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

Here, we have an hyperplane as a model:

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

In the case of \mathbb{R}^2

We have the following function:

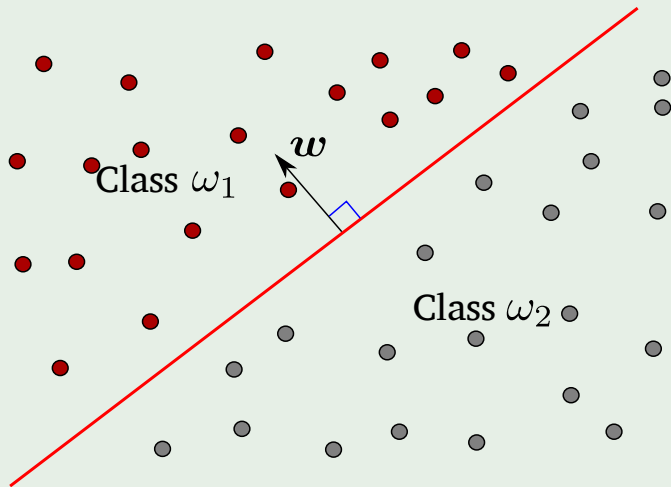
$$g(\mathbf{x}) = w_1 x_1 + w_2 x_2 + w_0 \quad (2)$$

Outline

- 1 Introduction
 - The Simplest Functions
 - **Splitting the Space**
 - The Decision Surface
 - Minimum Squared Error Procedure
 - The Error Idea
 - The Final Error Equation
 - The Data Matrix
 - Issues with Least Squares!!!

Splitting The Space \mathbb{R}^2

Using a simple straight line



Outline

1 Introduction

- The Simplest Functions
- Splitting the Space
- **The Decision Surface**
- Minimum Squared Error Procedure
- The Error Idea
- The Final Error Equation
- The Data Matrix
- Issues with Least Squares!!!

Defining a 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

Given 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

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

Defining a 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

Given 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

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

Defining a 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

Given 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

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

Defining a Decision Surface

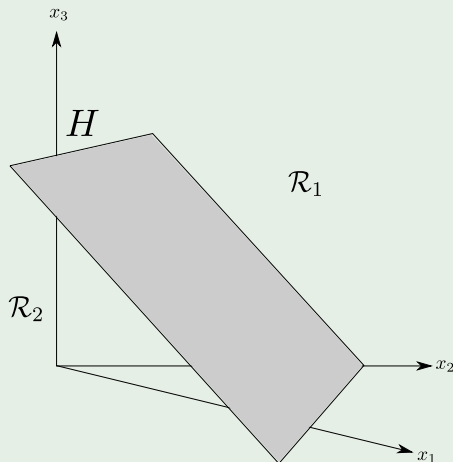
Thus

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

Remark: Any vector in the hyperplane is perpendicular to \mathbf{w}^T i.e. \mathbf{w}^T is normal to the hyperplane.

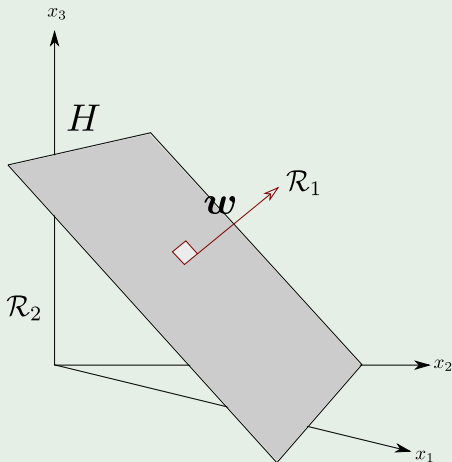
Therefore

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



Some Properties of the Hyperplane

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



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, you can give us a way to obtain the distance from x to the hyperplane H .

First, we express any x as follows

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

It is more

We can say the following

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

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

First, we express any x as follows

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

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

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

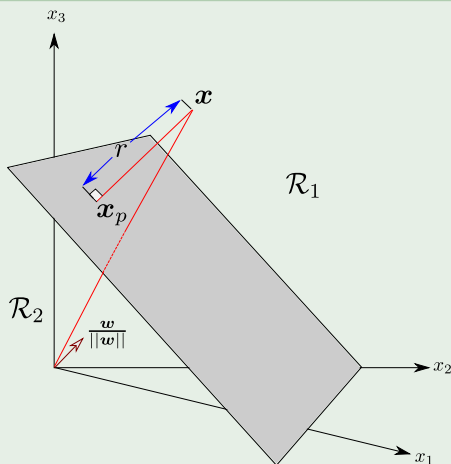
$$x = x_p + r \frac{w}{\|w\|}$$

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 $g(x_p) = 0$

We have that

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

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

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

(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.

In addition...

If we do the following

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

By making

$$x_0 = 1 \text{ and } \mathbf{y} = \begin{pmatrix} 1 \\ x_1 \\ \vdots \\ x_d \end{pmatrix} = \begin{pmatrix} 1 \\ \mathbf{x} \end{pmatrix}$$

Where

\mathbf{y} is called an augmented feature vector.

In addition...

If we do the following

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

By making

$$x_0 = 1 \text{ and } \mathbf{y} = \begin{pmatrix} 1 \\ x_1 \\ \vdots \\ x_d \end{pmatrix} = \begin{pmatrix} 1 \\ \mathbf{x} \end{pmatrix}$$

Where

\mathbf{y} is called an augmented feature vector.

In addition...

If we do the following

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

By making

$$x_0 = 1 \text{ and } \mathbf{y} = \begin{pmatrix} 1 \\ x_1 \\ \vdots \\ x_d \end{pmatrix} = \begin{pmatrix} 1 \\ \mathbf{x} \end{pmatrix}$$

Where

\mathbf{y} 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{y} 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{y} vectors, all lie in a d -dimensional subspace which is the \mathbf{x} -space itself.

Outline

- 1 Introduction
 - The Simplest Functions
 - Splitting the Space
 - The Decision Surface
 - **Minimum Squared Error Procedure**
 - The Error Idea
 - The Final Error Equation
 - The Data Matrix
 - Issues with Least Squares!!!

Initial Setup

Important

We get away from our initial normalization of the samples!!!

Now we are going to use the method know as

Minimum Squared Error

Initial Setup

Important

We get away from our initial normalization of the samples!!!

Now, we are going to use the method know as

Minimum Squared Error

Now, assume the following

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

① They are linearly separable!!!

② You require to label them.

Now, assume the following

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

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

We have a problem!!!

Which is the problem?

Now, assume the following

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

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

We have a problem!!!

Which is the problem?

We do not know the hyperplane!!!

Thus, what distance each point has to the hyperplane?

Now, assume the following

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

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

We have a problem!!!

Which is the problem?

We do not know the hyperplane!!!

Thus, what distance each point has to the hyperplane?

A Simple Solution For Our Quandary

Label the Classes

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

A Simple Solution For Our Quandary

Label the Classes

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

We produce the following labels

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

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

A Simple Solution For Our Quandary

Label the Classes

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

We produce the following labels

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

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

A Simple Solution For Our Quandary

Label the Classes

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

We produce the following labels

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

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

Outline

1 Introduction

- The Simplest Functions
- Splitting the Space
- The Decision Surface
- Minimum Squared Error Procedure
- **The Error Idea**
- The Final Error Equation
- The Data Matrix
- Issues with Least Squares!!!

Now, What?

Assume true function f is given by

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

Where the ϵ

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

Thus, we can do the following

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

Now, What?

Assume true function f is given by

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

Where the ϵ

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

Thus, we can do the following

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

Now, What?

Assume true function f is given by

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

Where the ϵ

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

Thus, we can do the following

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

Thus, we have

What to do?

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

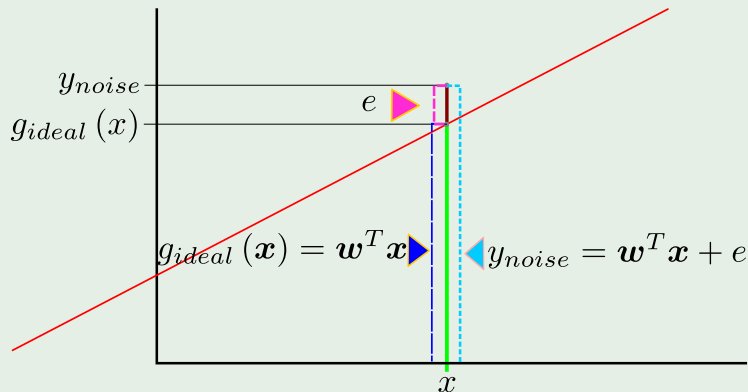
Graphically

Thus, we have

What to do?

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

Graphically



Outline

1 Introduction

- The Simplest Functions
- Splitting the Space
- The Decision Surface
- Minimum Squared Error Procedure
- The Error Idea
- **The Final Error Equation**
- The Data Matrix
- Issues with Least Squares!!!

Sum Over All Errors

We can do the following

$$J(\mathbf{w}) = \sum_{i=1}^N \epsilon_i^2 = \sum_{i=1}^N (y_i - g_{ideal}(\mathbf{x}))^2 \quad (11)$$

Remark: Known as least squares (Fitting the vertical offset!!!)

Sum Over All Errors

We can do the following

$$J(\mathbf{w}) = \sum_{i=1}^N \epsilon_i^2 = \sum_{i=1}^N (y_i - g_{ideal}(\mathbf{x}))^2 \quad (11)$$

Remark: Known as least squares (Fitting the vertical offset!!!)

Generalize

If

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

$$\sum_{i=1}^N (y_i - \mathbf{x}^T \mathbf{w})^2 = (\mathbf{y} - \mathbf{X}\mathbf{w})^T (\mathbf{y} - \mathbf{X}\mathbf{w}) = \|\mathbf{y} - \mathbf{X}\mathbf{w}\|_2^2 \quad (12)$$

Sum Over All Errors

We can do the following

$$J(\mathbf{w}) = \sum_{i=1}^N \epsilon_i^2 = \sum_{i=1}^N (y_i - g_{ideal}(\mathbf{x}))^2 \quad (11)$$

Remark: Known as least squares (Fitting the vertical offset!!!)

Generalize

If

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

• We have:

$$\sum_{i=1}^N (y_i - \mathbf{x}^T \mathbf{w})^2 = (\mathbf{y} - \mathbf{X}\mathbf{w})^T (\mathbf{y} - \mathbf{X}\mathbf{w}) = \|\mathbf{y} - \mathbf{X}\mathbf{w}\|_2^2 \quad (12)$$

Sum Over All Errors

We can do the following

$$J(\mathbf{w}) = \sum_{i=1}^N \epsilon_i^2 = \sum_{i=1}^N (y_i - g_{ideal}(\mathbf{x}))^2 \quad (11)$$

Remark: Known as least squares (Fitting the vertical offset!!!)

Generalize

If

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

$$\sum_{i=1}^N (y_i - \mathbf{x}^T \mathbf{w})^2 = (\mathbf{y} - \mathbf{X}\mathbf{w})^T (\mathbf{y} - \mathbf{X}\mathbf{w}) = \|\mathbf{y} - \mathbf{X}\mathbf{w}\|_2^2 \quad (12)$$

Outline

1 Introduction

- The Simplest Functions
- Splitting the Space
- The Decision Surface
- Minimum Squared Error Procedure
- The Error Idea
- The Final Error Equation
- **The Data Matrix**
- Issues with Least Squares!!!

What is \mathbf{X}

It is the Data Matrix

$$\mathbf{X} = \begin{pmatrix} 1 & (\mathbf{x}_1)_1 & \cdots & (\mathbf{x}_1)_j & \cdots & (\mathbf{x}_1)_d \\ \vdots & & & \vdots & & \vdots \\ 1 & (\mathbf{x}_i)_1 & & (\mathbf{x}_i)_j & & (\mathbf{x}_i)_d \\ \vdots & & & \vdots & & \vdots \\ 1 & (\mathbf{x}_N)_1 & \cdots & (\mathbf{x}_N)_j & \cdots & (\mathbf{x}_N)_d \end{pmatrix} \quad (13)$$

We know the following

$$\frac{d\mathbf{x}^T \mathbf{A} \mathbf{x}}{d\mathbf{x}} = \mathbf{A} \mathbf{x} + \mathbf{A}^T \mathbf{x}, \quad \frac{d\mathbf{A} \mathbf{x}}{d\mathbf{x}} = \mathbf{A} \quad (14)$$

What is \mathbf{X}

It is the Data Matrix

$$\mathbf{X} = \begin{pmatrix} 1 & (\mathbf{x}_1)_1 & \cdots & (\mathbf{x}_1)_j & \cdots & (\mathbf{x}_1)_d \\ \vdots & & & \vdots & & \vdots \\ 1 & (\mathbf{x}_i)_1 & & (\mathbf{x}_i)_j & & (\mathbf{x}_i)_d \\ \vdots & & & \vdots & & \vdots \\ 1 & (\mathbf{x}_N)_1 & \cdots & (\mathbf{x}_N)_j & \cdots & (\mathbf{x}_N)_d \end{pmatrix} \quad (13)$$

We know the following

$$\frac{d\mathbf{x}^T A \mathbf{x}}{d\mathbf{x}} = A\mathbf{x} + A^T \mathbf{x}, \quad \frac{dA\mathbf{x}}{d\mathbf{x}} = A \quad (14)$$

Note about other representations

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

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

We can expand our quadratic formula!!!

Thus

$$(y - Xw)^T (y - Xw) = y^T y - w^T X^T y - y^T X w + w^T X^T X w \quad (16)$$

Making Possible to have by deriving with respect to w and assuming that $X^T X$ is invertible

$$\hat{w} = (X^T X)^{-1} X^T y \quad (17)$$

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

Thus, How we get the discriminant function?

Any Ideas?

We can expand our quadratic formula!!!

Thus

$$(y - Xw)^T (y - Xw) = y^T y - w^T X^T y - y^T X w + w^T X^T X w \quad (16)$$

Making Possible to have by deriving with respect to w and assuming that $X^T X$ is invertible

$$\hat{w} = (X^T X)^{-1} X^T y \quad (17)$$

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

Thus, How we get the discriminant function?

Any Ideas?

We can expand our quadratic formula!!!

Thus

$$(y - Xw)^T (y - Xw) = y^T y - w^T X^T y - y^T X w + w^T X^T X w \quad (16)$$

Making Possible to have by deriving with respect to w and assuming that $X^T X$ is invertible

$$\hat{w} = (X^T X)^{-1} X^T y \quad (17)$$

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

Thus, How we get the discriminant function?

Any Ideas?

The Final Discriminant Function

Very Simple!!!

$$g(\mathbf{x}) = \mathbf{x}^T \hat{\mathbf{w}} = \mathbf{x}^T \left(\mathbf{X}^T \mathbf{X} \right)^{-1} \mathbf{X}^T \mathbf{y} \quad (18)$$

Pseudo-inverse of a Matrix

Definition

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

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

the pseudo inverse of A .

Reason

A^+ inverts A on its image

What?

If $w \in \text{image}(A)$, then there is some $v \in \mathbb{R}^n$ such that $w = Av$. Hence:

$$A^+ w = A^+ Av = (A^T A)^{-1} A^T Av$$

Pseudo-inverse of a Matrix

Definition

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

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

the pseudo inverse of A .

Reason

A^+ inverts A on its image

What?

If $w \in \text{image}(A)$, then there is some $v \in \mathbb{R}^n$ such that $w = Av$. Hence:

$$A^+ w = A^+ Av = (A^T A)^{-1} A^T Av$$

Pseudo-inverse of a Matrix

Definition

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

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

the pseudo inverse of A .

Reason

A^+ inverts A on its image

What?

If $w \in \text{image}(A)$, then there is some $v \in \mathbb{R}^n$ such that $w = Av$. Hence:

$$A^+ w = A^+ Av = (A^T A)^{-1} A^T Av$$

What lives where?

We have

- $\mathbf{X} \in \mathbb{R}^{N \times (d+1)}$
- $\text{Image}(\mathbf{X}) = \text{span}\{X_1^{\text{col}}, \dots, X_{d+1}^{\text{col}}\}$
- $x_i \in \mathbb{R}^d$
- $w \in \mathbb{R}^{d+1}$
- $X_i^{\text{col}}, y \in \mathbb{R}^N$

What lives where?

We have

- $\mathbf{X} \in \mathbb{R}^{N \times (d+1)}$
- $\text{Image}(\mathbf{X}) = \text{span} \{ \mathbf{X}_1^{\text{col}}, \dots, \mathbf{X}_{d+1}^{\text{col}} \}$
- $x_i \in \mathbb{R}^d$
- $w \in \mathbb{R}^{d+1}$
- $\mathbf{X}_i^{\text{col}}, y \in \mathbb{R}^N$

Basically, if the list of desired inputs the is being projected into

$$\text{span} \{ \mathbf{X}_1^{\text{col}}, \dots, \mathbf{X}_{d+1}^{\text{col}} \} \quad (19)$$

by the projection operator $\mathbf{X} (\mathbf{X}^T \mathbf{X})^{-1} \mathbf{X}^T$.

What lives where?

We have

- $\mathbf{X} \in \mathbb{R}^{N \times (d+1)}$
- $\text{Image}(\mathbf{X}) = \text{span} \{ \mathbf{X}_1^{\text{col}}, \dots, \mathbf{X}_{d+1}^{\text{col}} \}$
- $\mathbf{x}_i \in \mathbb{R}^d$
- $\mathbf{w} \in \mathbb{R}^{d+1}$
- $\mathbf{X}_i^{\text{col}}, \mathbf{y} \in \mathbb{R}^N$

Essentially, if the list of desired inputs the is being projected into

$$\text{span} \{ \mathbf{X}_1^{\text{col}}, \dots, \mathbf{X}_{d+1}^{\text{col}} \} \quad (19)$$

by the projection operator $\mathbf{X} (\mathbf{X}^T \mathbf{X})^{-1} \mathbf{X}^T$.

What lives where?

We have

- $\mathbf{X} \in \mathbb{R}^{N \times (d+1)}$
- $\text{Image}(\mathbf{X}) = \text{span} \{ \mathbf{X}_1^{\text{col}}, \dots, \mathbf{X}_{d+1}^{\text{col}} \}$
- $\mathbf{x}_i \in \mathbb{R}^d$
- $\mathbf{w} \in \mathbb{R}^{d+1}$
- $\mathbf{X}_i^{\text{col}}, \mathbf{y} \in \mathbb{R}^N$

Basically, \mathbf{y} , the list of desired inputs that is being protected into

$$\text{span} \{ \mathbf{X}_1^{\text{col}}, \dots, \mathbf{X}_{d+1}^{\text{col}} \} \quad (19)$$

by the projection operator $\mathbf{X} (\mathbf{X}^T \mathbf{X})^{-1} \mathbf{X}^T$.

What lives where?

We have

- $\mathbf{X} \in \mathbb{R}^{N \times (d+1)}$
- $\text{Image}(\mathbf{X}) = \text{span} \{ \mathbf{X}_1^{\text{col}}, \dots, \mathbf{X}_{d+1}^{\text{col}} \}$
- $\mathbf{x}_i \in \mathbb{R}^d$
- $\mathbf{w} \in \mathbb{R}^{d+1}$
- $\mathbf{X}_i^{\text{col}}, \mathbf{y} \in \mathbb{R}^N$

Essentially, \mathbf{y} , the list of desired inputs that is being protected into

$$\text{span} \{ \mathbf{X}_1^{\text{col}}, \dots, \mathbf{X}_{d+1}^{\text{col}} \} \quad (19)$$

by the projection operator $\mathbf{X} (\mathbf{X}^T \mathbf{X})^{-1} \mathbf{X}^T$.

What lives where?

We have

- $\mathbf{X} \in \mathbb{R}^{N \times (d+1)}$
- $\text{Image}(\mathbf{X}) = \text{span} \{ \mathbf{X}_1^{\text{col}}, \dots, \mathbf{X}_{d+1}^{\text{col}} \}$
- $\mathbf{x}_i \in \mathbb{R}^d$
- $\mathbf{w} \in \mathbb{R}^{d+1}$
- $\mathbf{X}_i^{\text{col}}, \mathbf{y} \in \mathbb{R}^N$

Basically \mathbf{y} , the list of desired inputs the is being protected into

$$\text{span} \{ \mathbf{X}_1^{\text{col}}, \dots, \mathbf{X}_{d+1}^{\text{col}} \} \quad (19)$$

by the projection operator $\mathbf{X} (\mathbf{X}^T \mathbf{X})^{-1} \mathbf{X}^T$.

Geometric Interpretation

We have

- 1 The image of the mapping w to Xw is a linear subspace of \mathbb{R}^N .
- 2 As w runs through all points \mathbb{R}^{d+1} , the function value Xw runs through all points in the image space
 $image(X) = span \{X_1^{col}, \dots, X_{d+1}^{col}\}$.
- 3 Each w defines one point $Xw = \sum_{j=0}^d w_j X_j^{col}$.
- 4 \hat{w} is the point which minimizes the distance $d(y, image(X))$.

Geometric Interpretation

We have

- 1 The image of the mapping w to Xw is a linear subspace of \mathbb{R}^N .
- 2 As w runs through all points \mathbb{R}^{d+1} , the function value Xw runs through all points in the image space
$$\text{image}(X) = \text{span}\{X_1^{\text{col}}, \dots, X_{d+1}^{\text{col}}\}.$$

Each w defines one point $Xw = \sum_{j=0}^d w_j X_j^{\text{col}}$.

\hat{w} is the point which minimizes the distance $d(y, \text{image}(X))$.

Geometric Interpretation

We have

- 1 The image of the mapping w to Xw is a linear subspace of \mathbb{R}^N .
 - 2 As w runs through all points \mathbb{R}^{d+1} , the function value Xw runs through all points in the image space
$$\text{image}(X) = \text{span} \{X_1^{\text{col}}, \dots, X_{d+1}^{\text{col}}\}.$$
 - 3 Each w defines one point $Xw = \sum_{j=0}^d w_j X_j^{\text{col}}$.
- ④ \hat{w} is the point which minimizes the distance $d(y, \text{image}(X))$.

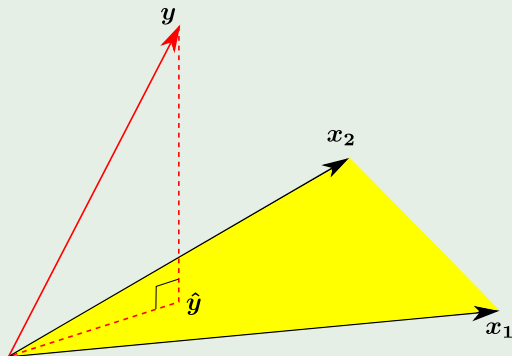
Geometric Interpretation

We have

- 1 The image of the mapping w to Xw is a linear subspace of \mathbb{R}^N .
- 2 As w runs through all points \mathbb{R}^{d+1} , the function value Xw runs through all points in the image space
$$\text{image}(X) = \text{span} \{X_1^{\text{col}}, \dots, X_{d+1}^{\text{col}}\}.$$
- 3 Each w defines one point $Xw = \sum_{j=0}^d w_j X_j^{\text{col}}$.
- 4 \hat{w} is the point which minimizes the distance $d(y, \text{image}(X))$.

Geometrically

Ahhhh!!!



Outline

- 1 Introduction
 - The Simplest Functions
 - Splitting the Space
 - The Decision Surface
 - Minimum Squared Error Procedure
 - The Error Idea
 - The Final Error Equation
 - The Data Matrix
 - **Issues with Least Squares!!!**

Issues with Least Squares

Robustness

- ① Least squares works only if \mathbf{X} has full column rank, i.e. if $\mathbf{X}^T \mathbf{X}$ is invertible.
- ② If $\mathbf{X}^T \mathbf{X}$ almost not invertible, least squares is numerically unstable.
 - ③ Statistical consequence: High variance of predictions.

Issues with Least Squares

Robustness

- 1 Least squares works only if X has full column rank, i.e. if $X^T X$ is invertible.
- 2 If $X^T X$ almost not invertible, least squares is numerically unstable.
 - Statistical consequence: High variance of predictions.

Not suited for high-dimensional data

- 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

Issues with Least Squares

Robustness

- 1 Least squares works only if X has full column rank, i.e. if $X^T X$ is invertible.
- 2 If $X^T X$ almost not invertible, least squares is numerically unstable.
 - 1 Statistical consequence: High variance of predictions.

Not suited for high-dimensional data

- 1 Modern problems: Many dimensions/features/predictors (possibly thousands).
- 2 Only a few of these may be important:
 - 1 It needs some form of feature selection.
 - 2 Possible some type of regularization

Why?

- 1 Treats all dimensions equally
- 2 Relevant dimensions are averaged with irrelevant ones

Issues with Least Squares

Robustness

- 1 Least squares works only if X has full column rank, i.e. if $X^T X$ is invertible.
- 2 If $X^T X$ almost not invertible, least squares is numerically unstable.
 - 1 Statistical consequence: High variance of predictions.

Not suited for high-dimensional data

- 1 Modern problems: Many dimensions/features/predictors (possibly thousands).
- 2 Only a few of these may be important:
 - 1 It needs some form of feature selection.
 - 2 Possible some type of regularization

Why?

- 2 Treats all dimensions equally
- 3 Relevant dimensions are averaged with irrelevant ones

Issues with Least Squares

Robustness

- 1 Least squares works only if X has full column rank, i.e. if $X^T X$ is invertible.
- 2 If $X^T X$ almost not invertible, least squares is numerically unstable.
 - 1 Statistical consequence: High variance of predictions.

Not suited for high-dimensional data

- 1 Modern problems: Many dimensions/features/predictors (possibly thousands).
- 2 Only a few of these may be important:
 - 1 It needs some form of feature selection.
 - 2 Possible some type of regularization

Why?

- 1 Treats all dimensions equally
- 2 Relevant dimensions are averaged with irrelevant ones

Issues with Least Squares

Robustness

- 1 Least squares works only if X has full column rank, i.e. if $X^T X$ is invertible.
- 2 If $X^T X$ almost not invertible, least squares is numerically unstable.
 - 1 Statistical consequence: High variance of predictions.

Not suited for high-dimensional data

- 1 Modern problems: Many dimensions/features/predictors (possibly thousands).
- 2 Only a few of these may be important:
 - 1 It needs some form of feature selection.

● Possible some type of regularization

Why?

- Treats all dimensions equally
- Relevant dimensions are averaged with irrelevant ones

Issues with Least Squares

Robustness

- 1 Least squares works only if X has full column rank, i.e. if $X^T X$ is invertible.
- 2 If $X^T X$ almost not invertible, least squares is numerically unstable.
 - 1 Statistical consequence: High variance of predictions.

Not suited for high-dimensional data

- 1 Modern problems: Many dimensions/features/predictors (possibly thousands).
- 2 Only a few of these may be important:
 - 1 It needs some form of feature selection.
 - 2 Possible some type of regularization

Why?

- Treats all dimensions equally
- Relevant dimensions are averaged with irrelevant ones

Issues with Least Squares

Robustness

- 1 Least squares works only if X has full column rank, i.e. if $X^T X$ is invertible.
- 2 If $X^T X$ almost not invertible, least squares is numerically unstable.
 - 1 Statistical consequence: High variance of predictions.

Not suited for high-dimensional data

- 1 Modern problems: Many dimensions/features/predictors (possibly thousands).
- 2 Only a few of these may be important:
 - 1 It needs some form of feature selection.
 - 2 Possible some type of regularization

Why?

- 1 Treats all dimensions equally
- 2 Relevant dimensions are averaged with irrelevant ones

Issues with Least Squares

Robustness

- 1 Least squares works only if X has full column rank, i.e. if $X^T X$ is invertible.
- 2 If $X^T X$ almost not invertible, least squares is numerically unstable.
 - 1 Statistical consequence: High variance of predictions.

Not suited for high-dimensional data

- 1 Modern problems: Many dimensions/features/predictors (possibly thousands).
- 2 Only a few of these may be important:
 - 1 It needs some form of feature selection.
 - 2 Possible some type of regularization

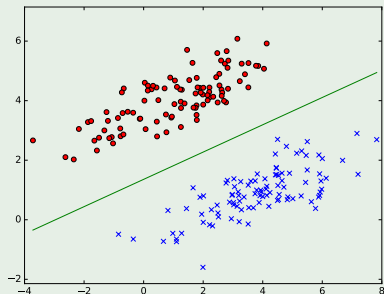
Why?

- 1 Treats all dimensions equally
- 2 Relevant dimensions are averaged with irrelevant ones

Issues with Least Squares

Problem with Outliers

No Outliers



Outliers

