

Introduction to Machine Learning

Regularization, Gradient Descent and Fisher Linear Discriminant

Andres Mendez-Vazquez

May 23, 2019

Outline

1 More in Regularization

- Introduction
- Smoothness of the Estimation
- The Error Estimate
- Choosing approximate inverses
- A Classic Example, Regularization as a Filter
- Another Example, The Landweber Iteration

2 Linear Regression using Gradient Descent

- Introduction
- What is the Gradient of the Equation?
- The Basic Algorithm
- How to obtain $\eta(k)$
 - Gold Section

3 The Gauss-Markov Theorem

- Statement
- Proof

4 Fisher Linear Discriminant

- History
- The Projection and The Rotation Idea
- Classifiers as Machines for dimensionality reduction
- Solution
 - Use the mean of each Class
 - Scatter measure
- The Cost Function
- A Transformation for simplification and defining the cost function
- Where is this used?
 - Applications
- Relation with Least Squared Error
 - What?

5 Exercises

- Some Stuff for you to try



Outline

1 More in Regularization

● Introduction

- Smoothness of the Estimation
- The Error Estimate
- Choosing approximate inverses
- A Classic Example, Regularization as a Filter
- Another Example, The Landweber Iteration

2 Linear Regression using Gradient Descent

- Introduction
- What is the Gradient of the Equation?
- The Basic Algorithm
- How to obtain $\eta(k)$
 - Gold Section

3 The Gauss-Markov Theorem

- Statement
- Proof

4 Fisher Linear Discriminant

- History
- The Projection and The Rotation Idea
- Classifiers as Machines for dimensionality reduction
- Solution
 - Use the mean of each Class
 - Scatter measure
- The Cost Function
- A Transformation for simplification and defining the cost function
- Where is this used?
 - Applications
- Relation with Least Squared Error
 - What?

5 Exercises

- Some Stuff for you to try



Well-Posed Problem

Definition by Hadamard (Circa 1902)

- Models of physical phenomenas should have the following properties
 - 1 A solution exists,
 - 2 The solution is unique,
 - 3 The solution's behavior changes continuously with the initial conditions.

Any other problem that fails in any of these conditions

- It is considered an Ill-Posed Problem.



Well-Posed Problem

Definition by Hadamard (Circa 1902)

- Models of physical phenomena should have the following properties
 - 1 A solution exists,
 - 2 The solution is unique,
 - 3 The solution's behavior changes continuously with the initial conditions.

Any other problem that fails in any of these conditions

- It is considered an Ill-Posed Problem.



Regularization in Linear Problems

In many applications of linear algebra

We want to find an estimation \hat{x} to a vector $x \in \mathbb{R}^d$ satisfying the approximation

$$Ax \approx y$$

When $A \in \mathbb{R}^{m \times d}$ is ill-conditioned or singular.

Regularization in Linear Problems

In many applications of linear algebra

We want to find an estimation \hat{x} to a vector $x \in \mathbb{R}^d$ satisfying the approximation

$$Ax \approx y$$

When $A \in \mathbb{R}^{m \times d}$ is ill-conditioned or singular.

The importance of the problem

The problems generating these situations are:

- Numerical differentiation of noisy data,
- Non parametric smoothing of curves and surfaces defined by scattered data,
- Image reconstruction,
- Inverse Laplace transforms,
- etc.

Regularization in Linear Problems

In many applications of linear algebra

We want to find an estimation \hat{x} to a vector $x \in \mathbb{R}^d$ satisfying the approximation

$$Ax \approx y$$

When $A \in \mathbb{R}^{m \times d}$ is ill-conditioned or singular.

The importance of the problem

The problems generating these situations are:

- Numerical differentiation of noisy data,
- Non parametric smoothing of curves and surfaces defined by scattered data,
- Image reconstruction,
- Inverse Laplace transforms,
- etc.

Regularization in Linear Problems

In many applications of linear algebra

We want to find an estimation \hat{x} to a vector $x \in \mathbb{R}^d$ satisfying the approximation

$$Ax \approx y$$

When $A \in \mathbb{R}^{m \times d}$ is ill-conditioned or singular.

The importance of the problem

The problems generating these situations are:

- Numerical differentiation of noisy data,
- Non parametric smoothing of curves and surfaces defined by scattered data,
- Image reconstruction,
- Inverse Laplace transforms,
- etc.

Regularization in Linear Problems

In many applications of linear algebra

We want to find an estimation \hat{x} to a vector $x \in \mathbb{R}^d$ satisfying the approximation

$$Ax \approx y$$

When $A \in \mathbb{R}^{m \times d}$ is ill-conditioned or singular.

The importance of the problem

The problems generating these situations are:

- 1 Numerical differentiation of noisy data,
- 2 Non parametric smoothing of curves and surfaces defined by scattered data,
- 3 Image reconstruction,
- 4 Inverse Laplace transforms,
- 5 etc.

Regularization in Linear Problems

In many applications of linear algebra

We want to find an estimation \hat{x} to a vector $x \in \mathbb{R}^d$ satisfying the approximation

$$Ax \approx y$$

When $A \in \mathbb{R}^{m \times d}$ is ill-conditioned or singular.

The importance of the problem

The problems generating these situations are:

- 1 Numerical differentiation of noisy data,
- 2 Non parametric smoothing of curves and surfaces defined by scattered data,
- 3 Image reconstruction,
- 4 Inverse Laplace transforms,
- 5 etc.

Regularization in Linear Problems

In many applications of linear algebra

We want to find an estimation \hat{x} to a vector $x \in \mathbb{R}^d$ satisfying the approximation

$$Ax \approx y$$

When $A \in \mathbb{R}^{m \times d}$ is ill-conditioned or singular.

The importance of the problem

The problems generating these situations are:

- 1 Numerical differentiation of noisy data,
- 2 Non parametric smoothing of curves and surfaces defined by scattered data,
- 3 Image reconstruction,
- 4 Inverse Laplace transforms,
- 5 etc.

Regularization in Linear Problems

In many applications of linear algebra

We want to find an estimation \hat{x} to a vector $x \in \mathbb{R}^d$ satisfying the approximation

$$Ax \approx y$$

When $A \in \mathbb{R}^{m \times d}$ is ill-conditioned or singular.

The importance of the problem

The problems generating these situations are:

- 1 Numerical differentiation of noisy data,
- 2 Non parametric smoothing of curves and surfaces defined by scattered data,
- 3 Image reconstruction,
- 4 Inverse Laplace transforms,

etc.

Regularization in Linear Problems

In many applications of linear algebra

We want to find an estimation \hat{x} to a vector $x \in \mathbb{R}^d$ satisfying the approximation

$$Ax \approx y$$

When $A \in \mathbb{R}^{m \times d}$ is ill-conditioned or singular.

The importance of the problem

The problems generating these situations are:

- 1 Numerical differentiation of noisy data,
- 2 Non parametric smoothing of curves and surfaces defined by scattered data,
- 3 Image reconstruction,
- 4 Inverse Laplace transforms,
- 5 etc.

In all such situations

The Vector \hat{x} generated by

① $\hat{x} = A^{-1}y$

② $\hat{x} = (A^T A)^{-1} A^T y$

if it exists at all

- It is usually a meaningless bad approximation to x .



In all such situations

The Vector \hat{x} generated by

① $\hat{x} = A^{-1}y$

② $\hat{x} = (A^T A)^{-1} A^T y$

If it exists at all

- It is usually a meaningless bad approximation to x .



Even

Even with an estimation $\hat{\mathbf{x}} = A\mathbf{y}$ as reasonable near to \mathbf{x}^* (Square Case)

$$\begin{aligned}\|\mathbf{x}^* - \hat{\mathbf{x}}\| &= \left\| A^{-1}A\mathbf{x}^* - A^{-1}\mathbf{y} \right\| \\ &\leq \|A^{-1}\| \|A\mathbf{x}^* - \mathbf{y}\| \quad \text{Holder's Inequality}\end{aligned}$$

- This Upper Bound is quite large.



Even

Even with an estimation $\hat{\mathbf{x}} = A\mathbf{y}$ as reasonable near to \mathbf{x}^* (Square Case)

$$\begin{aligned}\|\mathbf{x}^* - \hat{\mathbf{x}}\| &= \|A^{-1}A\mathbf{x}^* - A^{-1}\mathbf{y}\| \\ &\leq \|A^{-1}\| \|A\mathbf{x}^* - \mathbf{y}\| \quad \text{Holder's Inequality}\end{aligned}$$

- This Upper Bound is quite large.

With

- $\|A^{-1}\| = \sigma_{\max}(A)$ The largest singular value of matrix.

- $\|A\mathbf{x} - \mathbf{y}\| = \sqrt{(A\mathbf{x} - \mathbf{y})^T (A\mathbf{x} - \mathbf{y})}$

Even

Even with an estimation $\hat{\mathbf{x}} = A\mathbf{y}$ as reasonable near to \mathbf{x}^* (Square Case)

$$\begin{aligned}\|\mathbf{x}^* - \hat{\mathbf{x}}\| &= \|A^{-1}A\mathbf{x}^* - A^{-1}\mathbf{y}\| \\ &\leq \|A^{-1}\| \|A\mathbf{x}^* - \mathbf{y}\| \quad \text{Holder's Inequality}\end{aligned}$$

- This Upper Bound is quite large.

• $\|A^{-1}\| = \sigma_{\max}(A)$ The largest singular value of matrix.

• $\|A\mathbf{x} - \mathbf{y}\| = \sqrt{(A\mathbf{x} - \mathbf{y})^T (A\mathbf{x} - \mathbf{y})}$



Even

Even with an estimation $\hat{x} = A^{-1}y$ as reasonable near to x^* (Square Case)

$$\begin{aligned}\|x^* - \hat{x}\| &= \|A^{-1}Ax^* - A^{-1}y\| \\ &\leq \|A^{-1}\| \|Ax^* - y\| \quad \text{Holder's Inequality}\end{aligned}$$

- This Upper Bound is quite large.

With

- 1 $\|A^{-1}\| = \sigma_{\max}(A)$ The largest singular value of matrix.

$$\|Ax - y\| = \sqrt{(Ax - y)^T (Ax - y)}$$



Even

Even with an estimation $\hat{\mathbf{x}} = A\mathbf{y}$ as reasonable near to \mathbf{x}^* (Square Case)

$$\begin{aligned}\|\mathbf{x}^* - \hat{\mathbf{x}}\| &= \|A^{-1}A\mathbf{x}^* - A^{-1}\mathbf{y}\| \\ &\leq \|A^{-1}\| \|A\mathbf{x}^* - \mathbf{y}\| \quad \text{Holder's Inequality}\end{aligned}$$

- This Upper Bound is quite large.

With

- 1 $\|A^{-1}\| = \sigma_{\max}(A)$ The largest singular value of matrix.
- 2 $\|A\mathbf{x} - \mathbf{y}\| = \sqrt{(A\mathbf{x} - \mathbf{y})^T (A\mathbf{x} - \mathbf{y})}$

Therefore

Regularization techniques are needed to obtain meaningful solutions

- To problems that are called ill-posed problems.

Where some parameters are ill-determined

- By Least Square Methods
 - ▶ in particular when the number of parameters is larger than the number of available measurements!!!



Cinvestav

Therefore

Regularization techniques are needed to obtain meaningful solutions

- To problems that are called ill-posed problems.

Where some parameters are ill-determined

- By Least Square Methods
 - ▶ in particular when the number of parameters is larger than the number of available measurements!!!



Cinvestav

Outline

1 More in Regularization

- Introduction
- **Smoothness of the Estimation**
- The Error Estimate
- Choosing approximate inverses
- A Classic Example, Regularization as a Filter
- Another Example, The Landweber Iteration

2 Linear Regression using Gradient Descent

- Introduction
- What is the Gradient of the Equation?
- The Basic Algorithm
- How to obtain $\eta(k)$
 - Gold Section

3 The Gauss-Markov Theorem

- Statement
- Proof

4 Fisher Linear Discriminant

- History
- The Projection and The Rotation Idea
- Classifiers as Machines for dimensionality reduction
- Solution
 - Use the mean of each Class
 - Scatter measure
- The Cost Function
- A Transformation for simplification and defining the cost function
- Where is this used?
 - Applications
- Relation with Least Squared Error
 - What?

5 Exercises

- Some Stuff for you to try



Modeling Smoothness

Geometrically, regularization for smoothness means that

- We seek the least rough function that gives a certain degree of fit to the observed data.

Answer to this question is smoothness

- It is look at how many derivatives can be done before $\nabla^p f(x) = 0$



Modeling Smoothness

Geometrically, regularization for smoothness means that

- We seek the least rough function that gives a certain degree of fit to the observed data.

A way to measure smoothness

- It is look at how many derivatives can be done before $\nabla^p f(x) = 0$



Here, we want to model the idea of “Smoothness”

For this, we consider a continuous function f

- Where we use a vector w with features

$$w_i = f(t_i)$$

This, we can use a numerical differentiation method such that

$$w^{(1)} = \frac{df(t)}{dt}$$



Here, we want to model the idea of “Smoothness”

For this, we consider a continuous function f

- Where we use a vector w with features

$$w_i = f(t_i)$$

Thus, we can use a numerical differentiation method such that

$$w^{(1)} = \frac{df(t)}{dt}$$



Therefore, Assume Smoothness

We have a value such that $w = f(t)$

- Thus, we say that w is smooth “enough” if $w^{(1)} = \frac{df(t)}{dt}$ exists.

Now this can be repeated:

$$w^{(p)} = \frac{d^{(p)} f(t)}{dt^{(p)}}$$



Therefore, Assume Smoothness

We have a value such that $w = f(t)$

- Thus, we say that w is smooth “enough” if $w^{(1)} = \frac{df(t)}{dt}$ exists.

Now this can be repeated p

$$w^{(p)} = \frac{d^{(p)} f(t)}{dt^{(p)}}$$



Thus, it is possible to look at this smoothness

Using our Linear Algebra, we can represent this as a Linear Operator

$$w^{(p)} = Sw \text{ (The Smoothing Matrix)}$$



Thus

We can define the numerical differentiation of a $p + 1$ times

- Over a continuously differentiable function

$$y : [0, 1] \longrightarrow \mathbb{R}$$

Thus, finding our estimate $x(t) = y(t) = \nabla^p y(t)$

- Basically our problem of solving the linear system $Ax = y$

Or in other words

$$Ax(t) = \int_0^t x(\tau) d\tau$$



Thus

We can define the numerical differentiation of a $p + 1$ times

- Over a continuously differentiable function

$$y : [0, 1] \longrightarrow \mathbb{R}$$

Thus, finding our estimate $x(t) = y'(t) = \nabla y(t)$

- Basically our problem of solving the linear system $Ax = y$

Or in other words

$$Ax(t) = \int_0^t x(\tau) d\tau$$



Thus

We can define the numerical differentiation of a $p + 1$ times

- Over a continuously differentiable function

$$y : [0, 1] \longrightarrow \mathbb{R}$$

Thus, finding our estimate $x(t) = y'(t) = \nabla y(t)$

- Basically our problem of solving the linear system $Ax = y$

Or in other words

$$Ax(t) = \int_0^t x(\tau) d\tau$$



Therefore

The differentiability assumption says

$$\mathbf{w} = \nabla^{p+1}y = \nabla^p x \text{ is continuous and bounded}$$

Given that $\mathbf{w} = \nabla^p x$

- We may write the previous equation as

$$x = A^p \mathbf{w}$$



Therefore

The differentiability assumption says

$$\mathbf{w} = \nabla^{p+1}y = \nabla^p x \text{ is continuous and bounded}$$

Given that $A = \nabla^{-1}$

- We may write the previous equation as

$$x = A^p \mathbf{w}$$



Furthermore, Based in the following equalities

We can define the Adjoint Integral Operator is defined

$$\langle A^T x_1, x_2 \rangle = \langle x_1, Ax_2 \rangle$$

$$\langle x_1, x_2 \rangle = \int_0^1 x_1(t) x_2(t) dt$$

This with $y_1 = y$ and $y_2 = \nabla y$

$$\langle x_1, Ax_2 \rangle = \langle \nabla y_1, y_2 \rangle = -\langle y_1, \nabla y_2 \rangle = \langle -Ax_1, x_2 \rangle$$

- By Partial Integration

Then, under the following boundary conditions

- Assuming that y and its first $p+1$ derivatives vanish at $t=0$ and $t=1$.

Furthermore, Based in the following equalities

We can define the Adjoint Integral Operator is defined

$$\langle A^T x_1, x_2 \rangle = \langle x_1, Ax_2 \rangle$$

$$\langle x_1, x_2 \rangle = \int_0^1 x_1(t) x_2(t) dt$$

Thus with $Ax_i = y_i$ and $x_i = \nabla y_i$

$$\langle x_1, Ax_2 \rangle = \langle \nabla y_1, y_2 \rangle = -\langle y_1, \nabla y_2 \rangle = \langle -Ax_1, x_2 \rangle$$

- By Partial Integration

When under the following boundary conditions

- Assuming that y and its first $p+1$ derivatives vanish at $t=0$ and $t=1$.

Furthermore, Based in the following equalities

We can define the Adjoint Integral Operator is defined

$$\langle A^T x_1, x_2 \rangle = \langle x_1, Ax_2 \rangle$$

$$\langle x_1, x_2 \rangle = \int_0^1 x_1(t) x_2(t) dt$$

Thus with $Ax_i = y_i$ and $x_i = \nabla y_i$

$$\langle x_1, Ax_2 \rangle = \langle \nabla y_1, y_2 \rangle = -\langle y_1, \nabla y_2 \rangle = \langle -Ax_1, x_2 \rangle$$

- By Partial Integration

Then, under the following boundary conditions

- Assuming that y and its first $p + 1$ derivatives vanish at $t = 0$ and $t = 1$.

How?

We have

$$\langle \nabla y_1, y_2 \rangle = \int_0^1 \nabla y_1(t) y_2(t) dt$$

$$= y_1(t) y_2(t) \Big|_0^1 - \int_0^1 y_1(t) \nabla y_2(t) dt$$

$$= -\langle y_1, \nabla y_2 \rangle$$



How?

We have

$$\begin{aligned}\langle \nabla y_1, y_2 \rangle &= \int_0^1 \nabla y_1(t) y_2(t) dt \\ &= y_1(t) y_2(t) \Big|_0^1 - \int_0^1 y_1(t) \nabla y_2(t) dt \\ &= -\langle y_1, \nabla y_2 \rangle\end{aligned}$$

Then, if we assume that all entries in A are in \mathbb{R} ,

- $A^T = -A$



How?

We have

$$\begin{aligned}\langle \nabla y_1, y_2 \rangle &= \int_0^1 \nabla y_1(t) y_2(t) dt \\ &= y_1(t) y_2(t) \Big|_0^1 - \int_0^1 y_1(t) \nabla y_2(t) dt \\ &= -\langle y_1, \nabla y_2 \rangle\end{aligned}$$

Then, if we assume that all entries in A are in \mathbb{R}

- $A^T = -A$



Therefore

We have the following relation

$$\nabla y(t) = A^{-1}y(t)$$

Thus, it is possible to write the condition $y = Ay$ as $y = Sy$

- By absorbing the sign into w

$$S = \begin{cases} (A^T A)^{\frac{p}{2}} & \text{if } p \text{ is even} \\ (A^T A)^{\frac{p-1}{2}} A^T & \text{if } p \text{ is odd} \end{cases}$$

- For $p \geq 1$.



Therefore

We have the following relation

$$\nabla y(t) = A^{-1}y(t)$$

Thus, it is possible to write the condition $x = A^p w$ as $x = Sw$

- By absorbing the sign into w

$$S = \begin{cases} (A^T A)^{\frac{p}{2}} & \text{if } p \text{ is even} \\ (A^T A)^{\frac{p-1}{2}} A^T & \text{if } p \text{ is odd} \end{cases}$$

- For $p \geq 1$.



Outline

1 More in Regularization

- Introduction
- Smoothness of the Estimation
- **The Error Estimate**
- Choosing approximate inverses
- A Classic Example, Regularization as a Filter
- Another Example, The Landweber Iteration

2 Linear Regression using Gradient Descent

- Introduction
- What is the Gradient of the Equation?
- The Basic Algorithm
- How to obtain $\eta(k)$
 - Gold Section

3 The Gauss-Markov Theorem

- Statement
- Proof

4 Fisher Linear Discriminant

- History
- The Projection and The Rotation Idea
- Classifiers as Machines for dimensionality reduction
- Solution
 - Use the mean of each Class
 - Scatter measure
- The Cost Function
- A Transformation for simplification and defining the cost function
- Where is this used?
 - Applications
- Relation with Least Squared Error
 - What?

5 Exercises

- Some Stuff for you to try



The key to the treatment of ill-posed Linear Systems

It is a process called regularization that replaces A^{-1} by a family $C_h, h > 0$

- Of approximate inverses of A in such a way that, as $h \rightarrow 0$, the product $C_h A \rightarrow I$ in an appropriately restricted sense.
 - ▶ The parameter h is called the regularization parameter.



Therefore

It is usually possible to choose the C_h such that

- For a suitable exponent p (often $p = 1$ or 2), the constants
 - 1 $\gamma_1 = \sup_{h>0} h \|C_h\|.$
 - 2 $\gamma_2 = \sup_{h>0} h^{-p} \|(I - C_h A) S\|$

They are finite and of reasonable size.

- From this... we have...



Therefore

It is usually possible to choose the C_h such that

- For a suitable exponent p (often $p = 1$ or 2), the constants
 - 1 $\gamma_1 = \sup_{h>0} h \|C_h\|.$
 - 2 $\gamma_2 = \sup_{h>0} h^{-p} \|(I - C_h A) S\|$

They are finite and of reasonable size

- From this... we have...



The Following Theorem

Theorem

- Suppose $x = Sw$, and $\|Ax - y\| \leq \Delta \|w\|$ for some $\Delta > 0$.
- Then γ_1 and γ_2 implies

$$\|x - C_h y\| \leq \left[\gamma_1 \frac{\Delta}{h} + \gamma_2 h^p \right] \|w\|$$



The Following Theorem

Theorem

- Suppose $x = Sw$, and $\|Ax - y\| \leq \Delta \|w\|$ for some $\Delta > 0$.
- Then γ_1 and γ_2 implies

$$\|x - C_h y\| \leq \left[\gamma_1 \frac{\Delta}{h} + \gamma_2 h^p \right] \|w\|$$



The Following Theorem

Theorem

- Suppose $x = Sw$, and $\|Ax - y\| \leq \Delta \|w\|$ for some $\Delta > 0$.
- Then γ_1 and γ_2 implies

$$\|x - C_h y\| \leq \left[\gamma_1 \frac{\Delta}{h} + \gamma_2 h^p \right] \|w\|$$



For Example

For a well-posed data fitting problem

- One with a well-conditioned normal equation matrix $A^T A$

The least squares estimate

- It has an error of the order of Δ .

For example

- $C_h = (A^T A)^{-1} A^T = A^+ \implies h^{-1} = \|A^+\| = O(1)$ with $\gamma_1 = 1$ and $\gamma_2 = 0$ independent of p



For Example

For a well-posed data fitting problem

- One with a well-conditioned normal equation matrix $A^T A$

The least squares estimate

- It has an error of the order of Δ .

For example

- $C_h = (A^T A)^{-1} A^T = A^+ \implies h^{-1} = \|A^+\| = O(1)$ with $\gamma_1 = 1$ and $\gamma_2 = 0$ independent of p



For Example

For a well-posed data fitting problem

- One with a well-conditioned normal equation matrix $A^T A$

The least squares estimate

- It has an error of the order of Δ .

For example

- $C_h = (A^T A)^{-1} A^T = A^+ \implies h^{-1} = \|A^+\| = O(1)$ with $\gamma_1 = 1$ and $\gamma_2 = 0$ independent of p



Outline

1 More in Regularization

- Introduction
- Smoothness of the Estimation
- The Error Estimate
- **Choosing approximate inverses**
- A Classic Example, Regularization as a Filter
- Another Example, The Landweber Iteration

2 Linear Regression using Gradient Descent

- Introduction
- What is the Gradient of the Equation?
- The Basic Algorithm
- How to obtain $\eta(k)$
 - Gold Section

3 The Gauss-Markov Theorem

- Statement
- Proof

4 Fisher Linear Discriminant

- History
- The Projection and The Rotation Idea
- Classifiers as Machines for dimensionality reduction
- Solution
 - Use the mean of each Class
 - Scatter measure
- The Cost Function
- A Transformation for simplification and defining the cost function
- Where is this used?
 - Applications
- Relation with Least Squared Error
 - What?

5 Exercises

- Some Stuff for you to try

Now, we have

When A is rank deficient or becomes increasingly ill-conditioned

- We may improve the condition by modifying $A^T A$.

The simplest way to achieve this is by adding a small multiple of the identity:

- Since $A^T A$ is symmetric and positive semidefinite,

The matrix $A^T A + h^2 I$ has its eigenvalues

- They are in the interval $[h^2, h^2 + \|A\|^2]$



Now, we have

When A is rank deficient or becomes increasingly ill-conditioned

- We may improve the condition by modifying $A^T A$.

The simplest way to achieve this is by adding a small multiple of the identity

- Since $A^T A$ is symmetric and positive semidefinite.

The matrix $A^T A + h^2 I$ has its eigenvalues

- They are in the interval $[h^2, h^2 + \|A\|^2]$



Now, we have

When A is rank deficient or becomes increasingly ill-conditioned

- We may improve the condition by modifying $A^T A$.

The simplest way to achieve this is by adding a small multiple of the identity

- Since $A^T A$ is symmetric and positive semidefinite.

The matrix $A^T A + h^2 I$ has its eigenvalues

- They are in the interval $[h^2, h^2 + \|A\|^2]$



Here

The Condition Number of a Positive Definite Matrix Σ

$$\text{cond}(\Sigma) = \frac{\lambda_{\max}(\Sigma)}{\lambda_{\min}(\Sigma)}$$

- What happens

Which is related to the Maximum Likelihood of a Gaussian Distribution under a restriction

$$\begin{aligned} \max ML(\Sigma) \\ \text{s.t. } \text{cond}(\Sigma) \leq k \end{aligned}$$

- “Condition Number Regularized Covariance Estimation” by Won et. al

Here

The Condition Number of a Positive Definite Matrix Σ

$$\text{cond}(\Sigma) = \frac{\lambda_{\max}(\Sigma)}{\lambda_{\min}(\Sigma)}$$

- What happens

Which is related to the Maximum Likelihood of a Gaussian Distribution under a restriction

$$\begin{aligned} \max ML(\Sigma) \\ \text{s.t. } \text{cond}(\Sigma) \leq k \end{aligned}$$

- “Condition Number Regularized Covariance Estimation” by Won et. al

Here the Condition Number

It is

$$\text{cond} \left(A^T A + h^2 I \right) \leq \frac{h^2 + \|A\|^2}{h^2}$$

Therefore, we have from the previous slides

$$\hat{x} = \left(A^T A + h^2 I \right)^{-1} A^T y$$

Formula first derived by Tikhonov in 1963

- “Solution of incorrectly formulated problems and the regularization method,” Soviet Math. Dokl. 4 (1963), pp. 1035–1038.



Here the Condition Number

It is

$$\text{cond} \left(A^T A + h^2 I \right) \leq \frac{h^2 + \|A\|^2}{h^2}$$

Therefore, we have from the previous slides

$$\hat{x} = \left(A^T A + h^2 I \right)^{-1} A^T y$$

Example first derived by Tikhonov in 1963

- “Solution of incorrectly formulated problems and the regularization method,” Soviet Math. Dokl. 4 (1963), pp. 1035–1038.



Here the Condition Number

It is

$$\text{cond} \left(A^T A + h^2 I \right) \leq \frac{h^2 + \|A\|^2}{h^2}$$

Therefore, we have from the previous slides

$$\hat{x} = \left(A^T A + h^2 I \right)^{-1} A^T y$$

Formula first derived by Tikhonov in 1963

- “Solution of incorrectly formulated problems and the regularization method,” Soviet Math. Dokl. 4 (1963), pp. 1035–1038.



Finally

Corresponds to the family of approximate inverses (Tikhonov Regularization)

$$C_h = (A^T A + h^2 I)^{-1} A^T$$



Outline

1 More in Regularization

- Introduction
- Smoothness of the Estimation
- The Error Estimate
- Choosing approximate inverses
- **A Classic Example, Regularization as a Filter**
- Another Example, The Landweber Iteration

2 Linear Regression using Gradient Descent

- Introduction
- What is the Gradient of the Equation?
- The Basic Algorithm
- How to obtain $\eta(k)$
 - Gold Section

3 The Gauss-Markov Theorem

- Statement
- Proof

4 Fisher Linear Discriminant

- History
- The Projection and The Rotation Idea
- Classifiers as Machines for dimensionality reduction
- Solution
 - Use the mean of each Class
 - Scatter measure
- The Cost Function
- A Transformation for simplification and defining the cost function
- Where is this used?
 - Applications
- Relation with Least Squared Error
 - What?

5 Exercises

- Some Stuff for you to try



A Classic Example, The Finite-Dimensional Case

Given a Matrix K of $N \times N$

- with decomposition

$$K = Q\Sigma Q^t$$

- Such that $QQ^T = I$

Where

- Σ is the matrix $\text{diag}(\sigma_1, \sigma_2, \dots, \sigma_N)$ of eigenvalues with $\sigma_1 \geq \sigma_2 \geq \dots \geq \sigma_N$
- $Q = \begin{bmatrix} q_1 & q_2 & \dots & q_N \end{bmatrix}$ the corresponding eigenvectors.

Then, it is possible to write the following estimation

$$\hat{x} = K^{-1}Y = Q\Sigma^{-1}Q^T y = \sum_{i=1}^n \frac{1}{\sigma_i} \langle q_i, Y \rangle q_i$$

A Classic Example, The Finite-Dimensional Case

Given a Matrix K of $N \times N$

- with decomposition

$$K = Q\Sigma Q^t$$

- Such that $QQ^T = I$

Where

- Σ is the matrix $\text{diag}(\sigma_1, \sigma_2, \dots, \sigma_N)$ of eigenvalues with $\sigma_1 \geq \sigma_2 \geq \dots \geq \sigma_N$
- $Q = \begin{bmatrix} q_1 & q_2 & \cdots & q_N \end{bmatrix}$ the corresponding eigenvectors.

Then, it is possible to write the following estimation

$$\hat{x} = K^{-1}Y = Q\Sigma^{-1}Q^T y = \sum_{i=1}^n \frac{1}{\sigma_i} \langle q_i, Y \rangle q_i$$

A Classic Example, The Finite-Dimensional Case

Given a Matrix K of $N \times N$

- with decomposition

$$K = Q\Sigma Q^t$$

- Such that $QQ^T = I$

Where

- Σ is the matrix $diag(\sigma_1, \sigma_2, \dots, \sigma_N)$ of eigenvalues with $\sigma_1 \geq \sigma_2 \geq \dots \geq \sigma_N$
- $Q = \begin{bmatrix} q_1 & q_2 & \cdots & q_N \end{bmatrix}$ the corresponding eigenvectors.

Then, it is possible to write the following estimation

$$\hat{\mathbf{x}} = K^{-1}Y = Q\Sigma^{-1}Q^T\mathbf{y} = \sum_{i=1}^n \frac{1}{\sigma_i} \langle q_i, Y \rangle q_i$$

Therefore

If we start to see really small σ_i , the solution will be unstable

- It is more, if there are zero eigenvalues, the matrix will be impossible to invert.

Clearly, the coefficients of x will go infinity

$$x_i = \frac{1}{\sigma_i} (q_i, Y) \rightarrow \infty$$

- Or Statistical High Variance...



Therefore

If we start to see really small σ_i , the solution will be unstable

- It is more, if there are zero eigenvalues, the matrix will be impossible to invert.

Clearly, the coefficients of \hat{x} will go infinity

$$x_i = \frac{1}{\sigma_i} \langle q_i, Y \rangle \rightarrow \infty$$

- Or Statistical High Variance...



A Classic By Tikhonov

Add an extra term λ to avoid such problems

$$\hat{\mathbf{x}} = (K + n\lambda I)^{-1} Y = Q\Sigma^{-1}Q^T \mathbf{y}$$

Again simple linear algebra

- The eigenvalues are padded by the same value, and we do not care about the effect in the eigenvectors given that we care only in the directions!!!



A Classic By Tikhonov

Add an extra term λ to avoid such problems

$$\hat{\mathbf{x}} = (K + n\lambda I)^{-1} Y = Q\Sigma^{-1}Q^T \mathbf{y}$$

Again simple linear algebra

- The eigenvalues are padded by the same value, and we do not care about the effect in the eigenvectors given that we care only in the directions!!!



Thus

If we rewrite the equations

$$\hat{\mathbf{x}} = Q (\Sigma + n\lambda I)^{-1} Q^T \mathbf{y} = \sum_{i=1}^n \frac{1}{\sigma_i + n\lambda} \langle q_i, Y \rangle q_i$$

Actually, regularization filters out the undesired components

- If $\sigma_i \gg \lambda n$ then $\frac{1}{\sigma_i + n\lambda} \sim \frac{1}{\sigma_i}$
- If $\sigma_i \ll \lambda n$ then $\frac{1}{\sigma_i + n\lambda} \sim \frac{1}{n\lambda}$



Thus

If we rewrite the equations

$$\hat{\mathbf{x}} = Q (\Sigma + n\lambda I)^{-1} Q^T \mathbf{y} = \sum_{i=1}^n \frac{1}{\sigma_i + n\lambda} \langle q_i, Y \rangle q_i$$

Actually, regularization filters out the undesired components

- If $\sigma_i \gg \lambda n$ then $\frac{1}{\sigma_i + n\lambda} \sim \frac{1}{\sigma_i}$
- If $\sigma_i \ll \lambda n$ then $\frac{1}{\sigma_i + n\lambda} \sim \frac{1}{n\lambda}$



In a more general setup

Let be $G_\lambda(\sigma)$ a regularization function for the eigenvalues, we can then decompose K as

$$G_\lambda(K) = QG_\lambda(\sigma)Q^T$$

Therefore our estimation finishes as

$$G_\lambda(K)y = \sum_{i=1}^n G_\lambda(\sigma) \langle q_i, Y \rangle q_i$$



In a more general setup

Let be $G_\lambda(\sigma)$ a regularization function for the eigenvalues, we can then decompose K as

$$G_\lambda(K) = QG_\lambda(\sigma)Q^T$$

Therefore our estimation, finishes as

$$G_\lambda(K) \mathbf{y} = \sum_{i=1}^n G_\lambda(\sigma) \langle q_i, Y \rangle q_i$$



Clearly

For Tikhonov

$$G_{\lambda}(\sigma) = \frac{1}{\sigma_i + n\lambda}$$



Cinvestav

First

- In the inverse problems literature, many algorithms are known besides Tikhonov regularization.

These algorithms are defined by a suitable \mathcal{R} :

- They are not necessarily based on Regularized Empirical Risk Minimization (ERM):

$$R_{emp}(f) = \frac{1}{n} \sum_{i=1}^n L(f(x_i), y_i)$$

- However, they perform spectral regularization (Eigenvalue Based Regularization).



Remarks

First

- In the inverse problems literature, many algorithms are known besides Tikhonov regularization.

These algorithms are defined by a suitable G

- They are not necessarily based on **Regularized Empirical Risk Minimization (ERM)**:

$$R_{emp}(f) = \frac{1}{n} \sum_{i=1}^n L(f(x_i), y_i)$$

- However, they perform spectral regularization (Eigenvalue Based Regularization).



Spectral Filtering

Examples

- 1 Gradient Descent (or Landweber Iteration or L_2 Boosting)
- 2 ν -accelerated Landweber
- 3 Iterated Tikhonov Regularization
- 4 Truncated Singular Value Decomposition (TSVD)
- 5 Principle Component Regression (PCR)



Outline

1 More in Regularization

- Introduction
- Smoothness of the Estimation
- The Error Estimate
- Choosing approximate inverses
- A Classic Example, Regularization as a Filter
- **Another Example, The Landweber Iteration**

2 Linear Regression using Gradient Descent

- Introduction
- What is the Gradient of the Equation?
- The Basic Algorithm
- How to obtain $\eta(k)$
 - Gold Section

3 The Gauss-Markov Theorem

- Statement
- Proof

4 Fisher Linear Discriminant

- History
- The Projection and The Rotation Idea
- Classifiers as Machines for dimensionality reduction
- Solution
 - Use the mean of each Class
 - Scatter measure
- The Cost Function
- A Transformation for simplification and defining the cost function
- Where is this used?
 - Applications
- Relation with Least Squared Error
 - What?

5 Exercises

- Some Stuff for you to try

The Landweber Iteration

The Landweber iteration or Landweber algorithm

- It is an algorithm to solve ill-posed linear inverse problems

History

- The method was first proposed in the 1950s by Louis Landweber,

Remarks

- When A is nonsingular, then an explicit solution is $x = A^{-1}y$



The Landweber Iteration

The Landweber iteration or Landweber algorithm

- It is an algorithm to solve ill-posed linear inverse problems

It is quite old...

- The method was first proposed in the 1950s by Louis Landweber,

- When A is nonsingular, then an explicit solution is $x = A^{-1}y$



The Landweber Iteration

The Landweber iteration or Landweber algorithm

- It is an algorithm to solve ill-posed linear inverse problems

It is quite old...

- The method was first proposed in the 1950s by Louis Landweber,

Remarks

- When A is nonsingular, then an explicit solution is $x = A^{-1}y$



Therefore

The Landweber algorithm is an attempt to regularize the problem

- The algorithm tries to solve the minimization

$$\min_w \frac{\|y - Xw\|_2^2}{2}$$

Using the update

$$w_{k+1} = w_k + \eta X^T (y - Xw_k)$$

- where $0 < \eta < 2 \|X^T X\|_2^{-1} = 2\sigma$

This is given by the taking in account

$$\phi(w) = \frac{\|y - Xw\|_2^2}{2}$$

Therefore

The Landweber algorithm is an attempt to regularize the problem

- The algorithm tries to solve the minimization

$$\min_w \frac{\|\mathbf{y} - \mathbf{X}\mathbf{w}\|_2^2}{2}$$

Using the update

$$\mathbf{w}_{k+1} = \mathbf{w}_k + \eta \mathbf{X}^T (\mathbf{y} - \mathbf{X}\mathbf{w}_k)$$

- where $0 < \eta < 2 \left\| \mathbf{X}^T \mathbf{X} \right\|_2^{-1} = 2\sigma$

This is given by the taking in account

$$\phi(\mathbf{w}) = \frac{\|\mathbf{y} - \mathbf{X}\mathbf{w}\|_2^2}{2}$$

Therefore

The Landweber algorithm is an attempt to regularize the problem

- The algorithm tries to solve the minimization

$$\min_{\mathbf{w}} \frac{\|\mathbf{y} - \mathbf{X}\mathbf{w}\|_2^2}{2}$$

Using the update

$$\mathbf{w}_{k+1} = \mathbf{w}_k + \eta \mathbf{X}^T (\mathbf{y} - \mathbf{X}\mathbf{w}_k)$$

- where $0 < \eta < 2 \left\| \mathbf{X}^T \mathbf{X} \right\|_2^{-1} = 2\sigma$

This is given by the taking in account

$$\phi(\mathbf{w}) = \frac{\|\mathbf{y} - \mathbf{X}\mathbf{w}\|_2^2}{2}$$

Then

It is possible to show that the gradient of $\phi(\mathbf{w})$ is

$$\nabla \phi(\mathbf{w}) = -\mathbf{X}^T (\mathbf{y} - \mathbf{X}\mathbf{w}_k)$$

Therefore

- Each step in Landweber's method is a step in the direction of steepest descent.



Then

It is possible to show that the gradient of $\phi(\mathbf{w})$

$$\phi(\mathbf{w}) = -\mathbf{X}^T (\mathbf{y} - \mathbf{X}\mathbf{w}_k)$$

Therefore

- Each step in Landweber's method is a step in the direction of steepest descent.



Outline

1 More in Regularization

- Introduction
- Smoothness of the Estimation
- The Error Estimate
- Choosing approximate inverses
- A Classic Example, Regularization as a Filter
- Another Example, The Landweber Iteration

2 Linear Regression using Gradient Descent

- **Introduction**
- What is the Gradient of the Equation?
- The Basic Algorithm
- How to obtain $\eta(k)$
 - Gold Section

3 The Gauss-Markov Theorem

- Statement
- Proof

4 Fisher Linear Discriminant

- History
- The Projection and The Rotation Idea
- Classifiers as Machines for dimensionality reduction
- Solution
 - Use the mean of each Class
 - Scatter measure
- The Cost Function
- A Transformation for simplification and defining the cost function
- Where is this used?
 - Applications
- Relation with Least Squared Error
 - What?

5 Exercises

- Some Stuff for you to try

Given that the Canonical Solution has problems

We can develop a more robust algorithm

Using the Gradient Descent Idea

Specifically: The Gradient Descent

It uses the change in the surface of the cost function to obtain a direction of improvement.



Cinvestav

Given that the Canonical Solution has problems

We can develop a more robust algorithm

Using the Gradient Descent Idea

Basically, The Gradient Descent

It uses the change in the surface of the cost function to obtain a direction of improvement.



Gradient Descent

The basic procedure is as follow

- 1 Start with a random weight vector $w(1)$.
- 2 Compute the gradient vector $\nabla J(w(1))$.
- 3 Obtain value $w(2)$ by moving from $w(1)$ in the direction of the steepest descent:

$$w(k+1) = w(k) - \eta(k) \nabla J(w(k)) \quad (1)$$

$\eta(k)$ is a positive scale factor or learning rate!!!



Gradient Descent

The basic procedure is as follow

- 1 Start with a random weight vector $\boldsymbol{w}(1)$.
- 2 Compute the gradient vector $\nabla J(\boldsymbol{w}(1))$.
- 3 Obtain value $\boldsymbol{w}(2)$ by moving from $\boldsymbol{w}(1)$ in the direction of the steepest descent:

$$\boldsymbol{w}(k+1) = \boldsymbol{w}(k) - \eta(k) \nabla J(\boldsymbol{w}(k)) \quad (1)$$

$\eta(k)$ is a positive scale factor or learning rate!!!



Gradient Descent

The basic procedure is as follow

- 1 Start with a random weight vector $\boldsymbol{w}(1)$.
- 2 Compute the gradient vector $\nabla J(\boldsymbol{w}(1))$.
- 3 Obtain value $\boldsymbol{w}(2)$ by moving from $\boldsymbol{w}(1)$ in the direction of the steepest descent:

$$\boldsymbol{w}(k+1) = \boldsymbol{w}(k) - \eta(k) \nabla J(\boldsymbol{w}(k)) \quad (1)$$

$\eta(k)$ is a positive scale factor or learning rate!!!



Gradient Descent

The basic procedure is as follow

- 1 Start with a random weight vector $\mathbf{w}(1)$.
- 2 Compute the gradient vector $\nabla J(\mathbf{w}(1))$.
- 3 Obtain value $\mathbf{w}(2)$ by moving from $\mathbf{w}(1)$ in the direction of the steepest descent:

$$\mathbf{w}(k+1) = \mathbf{w}(k) - \eta(k) \nabla J(\mathbf{w}(k)) \quad (1)$$

$\eta(k)$ is a positive scale factor or learning rate!!!



Gradient Descent

The basic procedure is as follow

- 1 Start with a random weight vector $\mathbf{w}(1)$.
- 2 Compute the gradient vector $\nabla J(\mathbf{w}(1))$.
- 3 Obtain value $\mathbf{w}(2)$ by moving from $\mathbf{w}(1)$ in the direction of the steepest descent:

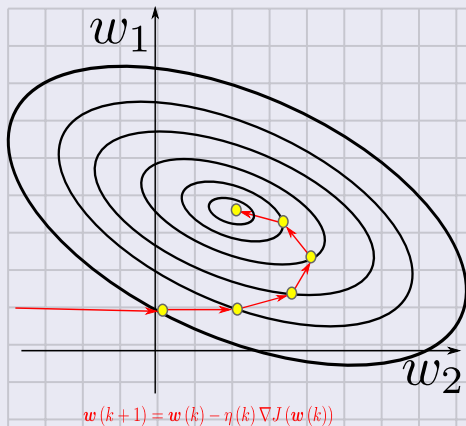
$$\mathbf{w}(k+1) = \mathbf{w}(k) - \eta(k) \nabla J(\mathbf{w}(k)) \quad (1)$$

$\eta(k)$ is a positive scale factor or learning rate!!!



Geometrically

We have the following



Outline

1 More in Regularization

- Introduction
- Smoothness of the Estimation
- The Error Estimate
- Choosing approximate inverses
- A Classic Example, Regularization as a Filter
- Another Example, The Landweber Iteration

2 Linear Regression using Gradient Descent

- Introduction
- **What is the Gradient of the Equation?**
- The Basic Algorithm
- How to obtain $\eta(k)$
 - Gold Section

3 The Gauss-Markov Theorem

- Statement
- Proof

4 Fisher Linear Discriminant

- History
- The Projection and The Rotation Idea
- Classifiers as Machines for dimensionality reduction
- Solution
 - Use the mean of each Class
 - Scatter measure
- The Cost Function
- A Transformation for simplification and defining the cost function
- Where is this used?
 - Applications
- Relation with Least Squared Error
 - What?

5 Exercises

- Some Stuff for you to try

For our full regularized equation

We have

$$J(\mathbf{w}) = \frac{1}{2} \sum_{i=1}^N \left(y_i - \sum_{j=1}^{d+1} x_j^i w_j \right)^2 + \frac{\lambda}{2} \sum_{j=1}^{d+1} w_j^2 \quad (2)$$

Then, for each w_j

$$\frac{dJ(\mathbf{w})}{dw_j} = - \sum_{i=1}^N \left[\left(y_i - \sum_{j=1}^{d+1} x_j^i w_j \right) x_j^i \right] + \lambda w_j \quad (3)$$

Therefore

$$\nabla J(\mathbf{w}(k)) = \begin{pmatrix} - \sum_{i=1}^N \left[\left(y_i - \sum_{j=1}^{d+1} x_j^i w_j \right) x_1^i \right] + \lambda w_1 \\ \vdots \\ - \sum_{i=1}^N \left[\left(y_i - \sum_{j=1}^{d+1} x_j^i w_j \right) x_{d+1}^i \right] + \lambda w_{d+1} \end{pmatrix}$$

For our full regularized equation

We have

$$J(\mathbf{w}) = \frac{1}{2} \sum_{i=1}^N \left(y_i - \sum_{j=1}^{d+1} x_j^i w_j \right)^2 + \frac{\lambda}{2} \sum_{j=1}^{d+1} w_j^2 \quad (2)$$

Then, for each w_j

$$\frac{dJ(\mathbf{w})}{dw_j} = - \sum_{i=1}^N \left[\left(y_i - \sum_{j=1}^{d+1} x_j^i w_j \right) x_j^i \right] + \lambda w_j \quad (3)$$

Therefore

$$\nabla J(\mathbf{w}(k)) = \begin{pmatrix} -\sum_{i=1}^N \left[\left(y_i - \sum_{j=1}^{d+1} x_j^i w_j \right) x_1^i \right] + \lambda w_1 \\ \vdots \\ -\sum_{i=1}^N \left[\left(y_i - \sum_{j=1}^{d+1} x_j^i w_j \right) x_{d+1}^i \right] + \lambda w_{d+1} \end{pmatrix}$$

For our full regularized equation

We have

$$J(\mathbf{w}) = \frac{1}{2} \sum_{i=1}^N \left(y_i - \sum_{j=1}^{d+1} x_j^i w_j \right)^2 + \frac{\lambda}{2} \sum_{j=1}^{d+1} w_j^2 \quad (2)$$

Then, for each w_j

$$\frac{dJ(\mathbf{w})}{dw_j} = - \sum_{i=1}^N \left[\left(y_i - \sum_{j=1}^{d+1} x_j^i w_j \right) x_j^i \right] + \lambda w_j \quad (3)$$

Therefore

$$\nabla J(\mathbf{w}(k)) = \begin{pmatrix} - \sum_{i=1}^N \left[\left(y_i - \sum_{j=1}^{d+1} x_j^i w_j \right) x_1^i \right] + \lambda w_1 \\ \vdots \\ - \sum_{i=1}^N \left[\left(y_i - \sum_{j=1}^{d+1} x_j^i w_j \right) x_{d+1}^i \right] + \lambda w_{d+1} \end{pmatrix}$$

Outline

- 1 More in Regularization
 - Introduction
 - Smoothness of the Estimation
 - The Error Estimate
 - Choosing approximate inverses
 - A Classic Example, Regularization as a Filter
 - Another Example, The Landweber Iteration
- 2 Linear Regression using Gradient Descent
 - Introduction
 - What is the Gradient of the Equation?
 - **The Basic Algorithm**
 - How to obtain $\eta(k)$
 - Gold Section
- 3 The Gauss-Markov Theorem
 - Statement
 - Proof
- 4 Fisher Linear Discriminant
 - History
 - The Projection and The Rotation Idea
 - Classifiers as Machines for dimensionality reduction
 - Solution
 - Use the mean of each Class
 - Scatter measure
 - The Cost Function
 - A Transformation for simplification and defining the cost function
 - Where is this used?
 - Applications
 - Relation with Least Squared Error
 - What?
- 5 Exercises
 - Some Stuff for you to try

Algorithm

Gradient Decent

- 1 Initialize w , criterion θ , $\eta(\cdot)$, $k = 0$
- 2 do $k = k + 1$
- 3 $w(k) = w(k-1) - \eta(k) \nabla J(w(k-1))$
- 4 until $\eta(k) \nabla J(w(k)) < \theta$
- 5 return w



Algorithm

Gradient Decent

- 1 Initialize w , criterion θ , $\eta(\cdot)$, $k = 0$
- 2 do $k = k + 1$
 - 3 $w(k) = w(k-1) - \eta(k) \nabla J(w(k-1))$
 - 4 until $\eta(k) \nabla J(w(k)) < \theta$
 - 5 return w

Problem!!! How to choose the learning rate?

- If $\eta(k)$ is too small, convergence is quite slow!!!
- If $\eta(k)$ is too large, correction will overshoot and can even diverge!!!



Algorithm

Gradient Decent

- 1 Initialize \mathbf{w} , criterion θ , $\eta(\cdot)$, $k = 0$
- 2 do $k = k + 1$
- 3 $\mathbf{w}(k) = \mathbf{w}(k - 1) - \eta(k) \nabla J(\mathbf{w}(k - 1))$
- 4 until $\eta(k) \nabla J(\mathbf{w}(k)) < \theta$
- 5 return \mathbf{w}

Problem!!! How to choose the learning rate?

- If $\eta(k)$ is too small, convergence is quite slow!!!
- If $\eta(k)$ is too large, correction will overshoot and can even diverge!!!



Algorithm

Gradient Decent

- 1 Initialize \mathbf{w} , criterion θ , $\eta(\cdot)$, $k = 0$
 - 2 do $k = k + 1$
 - 3 $\mathbf{w}(k) = \mathbf{w}(k - 1) - \eta(k) \nabla J(\mathbf{w}(k - 1))$
 - 4 until $\eta(k) \nabla J(\mathbf{w}(k)) < \theta$
- return \mathbf{w}

Problem!!! How to choose the learning rate?

- If $\eta(k)$ is too small, convergence is quite slow!!!
- If $\eta(k)$ is too large, correction will overshoot and can even diverge!!!



Algorithm

Gradient Decent

- 1 Initialize \mathbf{w} , criterion θ , $\eta(\cdot)$, $k = 0$
- 2 do $k = k + 1$
- 3 $\mathbf{w}(k) = \mathbf{w}(k - 1) - \eta(k) \nabla J(\mathbf{w}(k - 1))$
- 4 until $\eta(k) \nabla J(\mathbf{w}(k)) < \theta$
- 5 return \mathbf{w}

Problem!!! How to choose the learning rate?

- If $\eta(k)$ is too small, convergence is quite slow!!!
- If $\eta(k)$ is too large, correction will overshoot and can even diverge!!!



Algorithm

Gradient Decent

- 1 Initialize \mathbf{w} , criterion θ , $\eta(\cdot)$, $k = 0$
- 2 do $k = k + 1$
- 3 $\mathbf{w}(k) = \mathbf{w}(k - 1) - \eta(k) \nabla J(\mathbf{w}(k - 1))$
- 4 until $\eta(k) \nabla J(\mathbf{w}(k)) < \theta$
- 5 return \mathbf{w}

Problem!!! How to choose the learning rate?

- If $\eta(k)$ is too small, convergence is quite slow!!!
- If $\eta(k)$ is too large, correction will overshoot and can even diverge!!!



Algorithm

Gradient Decent

- 1 Initialize \mathbf{w} , criterion θ , $\eta(\cdot)$, $k = 0$
- 2 do $k = k + 1$
- 3 $\mathbf{w}(k) = \mathbf{w}(k - 1) - \eta(k) \nabla J(\mathbf{w}(k - 1))$
- 4 until $\eta(k) \nabla J(\mathbf{w}(k)) < \theta$
- 5 return \mathbf{w}

Problem!!! How to choose the learning rate?

- If $\eta(k)$ is too small, convergence is quite slow!!!
- If $\eta(k)$ is too large, correction will overshoot and can even diverge!!!



Outline

1 More in Regularization

- Introduction
- Smoothness of the Estimation
- The Error Estimate
- Choosing approximate inverses
- A Classic Example, Regularization as a Filter
- Another Example, The Landweber Iteration

2 Linear Regression using Gradient Descent

- Introduction
- What is the Gradient of the Equation?
- The Basic Algorithm
- **How to obtain $\eta(k)$**
 - Gold Section

3 The Gauss-Markov Theorem

- Statement
- Proof

4 Fisher Linear Discriminant

- History
- The Projection and The Rotation Idea
- Classifiers as Machines for dimensionality reduction
- Solution
 - Use the mean of each Class
 - Scatter measure
- The Cost Function
- A Transformation for simplification and defining the cost function
- Where is this used?
 - Applications
- Relation with Least Squared Error
 - What?

5 Exercises

- Some Stuff for you to try

Using the Taylor's second-order expansion around value $\mathbf{w}(k)$

We do the following

$$J(\mathbf{w}) = J(\mathbf{w}(k)) + \nabla J^T(\mathbf{w} - \mathbf{w}(k)) + \frac{1}{2}(\mathbf{w} - \mathbf{w}(k))^T \mathbf{H}(\mathbf{w} - \mathbf{w}(k)) \quad (4)$$

Remark: This is known as Taylor's Second Order expansion!!!



Using the Taylor's second-order expansion around value $w(k)$

We do the following

$$J(w) = J(w(k)) + \nabla J^T(w - w(k)) + \frac{1}{2}(w - w(k))^T H(w - w(k)) \quad (4)$$

Remark: This is known as Taylor's Second Order expansion!!!

- ∇J is the vector of partial derivatives $\frac{\partial J}{\partial w_i}$ evaluated at $w(k)$.
- H is the Hessian matrix of second partial derivatives $\frac{\partial^2 J}{\partial w_i \partial w_j}$ evaluated at $w(k)$.



Using the Taylor's second-order expansion around value $\mathbf{w}(k)$

We do the following

$$J(\mathbf{w}) = J(\mathbf{w}(k)) + \nabla J^T(\mathbf{w} - \mathbf{w}(k)) + \frac{1}{2}(\mathbf{w} - \mathbf{w}(k))^T \mathbf{H}(\mathbf{w} - \mathbf{w}(k)) \quad (4)$$

Remark: This is known as Taylor's Second Order expansion!!!

Here, we have

- ∇J is the vector of partial derivatives $\frac{\partial J}{\partial w_i}$ evaluated at $\mathbf{w}(k)$.
- \mathbf{H} is the Hessian matrix of second partial derivatives $\frac{\partial^2 J}{\partial w_i \partial w_j}$ evaluated at $\mathbf{w}(k)$.



Using the Taylor's second-order expansion around value $\mathbf{w}(k)$

We do the following

$$J(\mathbf{w}) = J(\mathbf{w}(k)) + \nabla J^T(\mathbf{w} - \mathbf{w}(k)) + \frac{1}{2}(\mathbf{w} - \mathbf{w}(k))^T \mathbf{H}(\mathbf{w} - \mathbf{w}(k)) \quad (4)$$

Remark: This is known as Taylor's Second Order expansion!!!

Here, we have

- ∇J is the vector of partial derivatives $\frac{\partial J}{\partial w_i}$ evaluated at $\mathbf{w}(k)$.
- \mathbf{H} is the Hessian matrix of second partial derivatives $\frac{\partial^2 J}{\partial w_i \partial w_j}$ evaluated at $\mathbf{w}(k)$.



Then

We substitute (Eq. 1) into (Eq. 4)

$$\mathbf{w}(k+1) - \mathbf{w}(k) = \eta(k) \nabla J(\mathbf{w}(k)) \quad (5)$$

We have then

$$J(\mathbf{w}(k+1)) \cong J(\mathbf{w}(k)) + \nabla J^T(-\eta(k) \nabla J(\mathbf{w}(k))) + \dots \\ + \frac{1}{2} (-\eta(k) \nabla J(\mathbf{w}(k)))^T H (-\eta(k) \nabla J(\mathbf{w}(k)))$$

Finally, we have

$$J(\mathbf{w}(k+1)) \cong J(\mathbf{w}(k)) - \eta(k) \|\nabla J\|^2 + \frac{1}{2} \eta^2(k) \nabla J^T H \nabla J \quad (6)$$



Then

We substitute (Eq. 1) into (Eq. 4)

$$\mathbf{w}(k+1) - \mathbf{w}(k) = \eta(k) \nabla J(\mathbf{w}(k)) \quad (5)$$

We have then

$$J(\mathbf{w}(k+1)) \cong J(\mathbf{w}(k)) + \nabla J^T(-\eta(k) \nabla J(\mathbf{w}(k))) + \dots \\ \frac{1}{2} (-\eta(k) \nabla J(\mathbf{w}(k)))^T \mathbf{H}(-\eta(k) \nabla J(\mathbf{w}(k)))$$

Finally, we have

$$J(\mathbf{w}(k+1)) \cong J(\mathbf{w}(k)) - \eta(k) \|\nabla J\|^2 + \frac{1}{2} \eta^2(k) \nabla J^T \mathbf{H} \nabla J \quad (6)$$



Then

We substitute (Eq. 1) into (Eq. 4)

$$\mathbf{w}(k+1) - \mathbf{w}(k) = \eta(k) \nabla J(\mathbf{w}(k)) \quad (5)$$

We have then

$$J(\mathbf{w}(k+1)) \cong J(\mathbf{w}(k)) + \nabla J^T(-\eta(k) \nabla J(\mathbf{w}(k))) + \dots \\ \frac{1}{2} (-\eta(k) \nabla J(\mathbf{w}(k)))^T \mathbf{H} (-\eta(k) \nabla J(\mathbf{w}(k)))$$

Finally, we have

$$J(\mathbf{w}(k+1)) \cong J(\mathbf{w}(k)) - \eta(k) \|\nabla J\|^2 + \frac{1}{2} \eta^2(k) \nabla J^T \mathbf{H} \nabla J \quad (6)$$



Derive with respect to $\eta(k)$ and make the result equal to zero

We have then

$$-\|\nabla J\|^2 + \eta(k) \nabla J^T \mathbf{H} \nabla J = 0 \quad (7)$$

Finally

$$\eta(k) = \frac{\|\nabla J\|^2}{\nabla J^T \mathbf{H} \nabla J} \quad (8)$$

Remark: This is the optimal step size!!!

Problem!!!

Calculating \mathbf{H} can be quite expensive!!!



Derive with respect to $\eta(k)$ and make the result equal to zero

We have then

$$-\|\nabla J\|^2 + \eta(k) \nabla J^T \mathbf{H} \nabla J = 0 \quad (7)$$

Finally

$$\eta(k) = \frac{\|\nabla J\|^2}{\nabla J^T \mathbf{H} \nabla J} \quad (8)$$

Remark This is the optimal step size!!!

Calculating \mathbf{H} can be quite expensive!!!



Derive with respect to $\eta(k)$ and make the result equal to zero

We have then

$$-\|\nabla J\|^2 + \eta(k) \nabla J^T \mathbf{H} \nabla J = 0 \quad (7)$$

Finally

$$\eta(k) = \frac{\|\nabla J\|^2}{\nabla J^T \mathbf{H} \nabla J} \quad (8)$$

Remark This is the optimal step size!!!

Problem!!!

Calculating \mathbf{H} can be quite expensive!!!



We can have an adaptive linear search!!!

We can use the idea of having everything fixed, but $\eta(k)$

Then, we can have the following function

$$f(\eta(k)) = J(\mathbf{w}(k) - \eta(k) \nabla J(\mathbf{w}(k)))$$

- We can optimized using linear search methods



We can have an adaptive linear search!!!

We can use the idea of having everything fixed, but $\eta(k)$

Then, we can have the following function

$$f(\eta(k)) = J(\mathbf{w}(k) - \eta(k) \nabla J(\mathbf{w}(k)))$$

- We can optimized using linear search methods

Linear Search Methods

- Backtracking linear search
- Bisection method
- Golden ratio
- Etc.



We can have an adaptive linear search!!!

We can use the idea of having everything fixed, but $\eta(k)$

Then, we can have the following function

$$f(\eta(k)) = J(\mathbf{w}(k) - \eta(k) \nabla J(\mathbf{w}(k)))$$

- We can optimized using linear search methods

Linear Search Methods

- Backtracking linear search
- Bisection method
- Golden ratio
- Etc.



We can have an adaptive linear search!!!

We can use the idea of having everything fixed, but $\eta(k)$

Then, we can have the following function

$$f(\eta(k)) = J(\mathbf{w}(k) - \eta(k) \nabla J(\mathbf{w}(k)))$$

- We can optimized using linear search methods

Linear Search Methods

- Backtracking linear search
- Bisection method
- Golden ratio
- Etc.



We can have an adaptive linear search!!!

We can use the idea of having everything fixed, but $\eta(k)$

Then, we can have the following function

$$f(\eta(k)) = J(\mathbf{w}(k) - \eta(k) \nabla J(\mathbf{w}(k)))$$

- We can optimized using linear search methods

Linear Search Methods

- Backtracking linear search
- Bisection method
- Golden ratio

• Etc.



We can have an adaptive linear search!!!

We can use the idea of having everything fixed, but $\eta(k)$

Then, we can have the following function

$$f(\eta(k)) = J(\mathbf{w}(k) - \eta(k) \nabla J(\mathbf{w}(k)))$$

- We can optimized using linear search methods

Linear Search Methods

- Backtracking linear search
- Bisection method
- Golden ratio
- Etc.



Outline

1 More in Regularization

- Introduction
- Smoothness of the Estimation
- The Error Estimate
- Choosing approximate inverses
- A Classic Example, Regularization as a Filter
- Another Example, The Landweber Iteration

2 Linear Regression using Gradient Descent

- Introduction
- What is the Gradient of the Equation?
- The Basic Algorithm
- **How to obtain $\eta(k)$**
 - **Gold Section**

3 The Gauss-Markov Theorem

- Statement
- Proof

4 Fisher Linear Discriminant

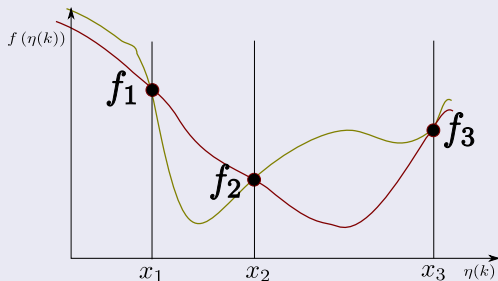
- History
- The Projection and The Rotation Idea
- Classifiers as Machines for dimensionality reduction
- Solution
 - Use the mean of each Class
 - Scatter measure
- The Cost Function
- A Transformation for simplification and defining the cost function
- Where is this used?
 - Applications
- Relation with Least Squared Error
 - What?

5 Exercises

- Some Stuff for you to try

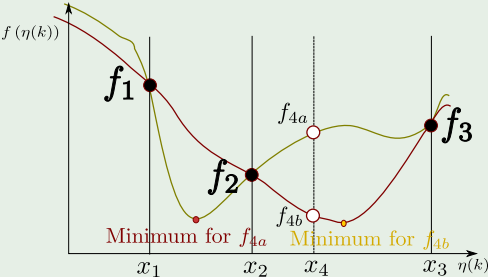
Gold Section

We have $f(\eta(k)) = J(\mathbf{w}(k) - \eta(k) \nabla J(\mathbf{w}(k)))$



Golden Section

Thus the idea is to use an evaluation f_4 to decide which subsection to drop



What is the Golden Ratio Idea?

Basically, given an interval $[x_1, x_3]$

Then, we select a point x_2 and x_3 such that we have a two possible intervals of search for the minimum

① $[x_1, x_4]$

② $[x_2, x_3]$

The Golden Linear Search requires these intervals be equal!!!

If they are not,

- You could run to a series of search wider intervals slowing down the rate of convergence.



What is the Golden Ratio Idea?

Basically, given an interval $[x_1, x_3]$

Then, we select a point x_2 and x_3 such that we have a two possible intervals of search for the minimum

① $[x_1, x_4]$

② $[x_2, x_3]$

The Golden Linear Search requires these intervals be equal!!!

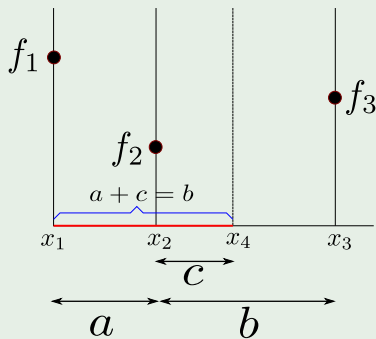
If they are not,

- You could run to a series of search wider intervals slowing down the rate of convergence.



How?

By the equality $b = a + c$



Therefore

We have the following question?

Where do you place x_2 ? Thus you can generate x_4

You want to avoid

- x_2 too close to x_1 or x_3



Cinvestav

Therefore

We have the following question?

Where do you place x_2 ? Thus you can generate x_4

You want to avoid

- x_2 to close to x_1 or x_3



The process is as follow

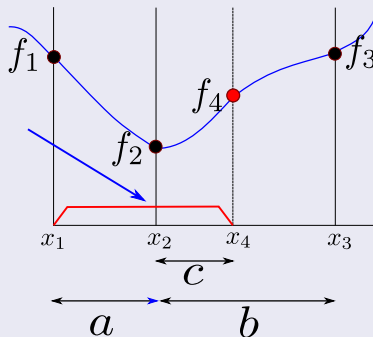
We define

- $f_1 = f(x_1)$
- $f_2 = f(x_2)$
- $f_3 = f(x_3)$
- $f_4 = f(x_4)$



Two Cases

If $f_2 < f_4$ then the minimum lies between x_1 and x_4 and the new triplet is x_1, x_2 and x_4 .



Here, we have the realization that

We have interval size reduction

$$x_4 - x_1 = \varphi(x_3 - x_1) \mapsto x_4 = x_1 + \varphi x_3 - \varphi x_1$$

Then

$$x_4 = (1 - \varphi)x_1 + \varphi x_3$$



Here, we have the realization that

We have interval size reduction

$$x_4 - x_1 = \varphi (x_3 - x_1) \mapsto x_4 = x_1 + \varphi x_3 - \varphi x_1$$

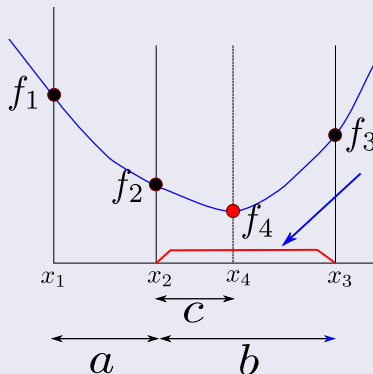
Then

$$x_4 = (1 - \varphi) x_1 + \varphi x_3$$



Two Cases

If $f_4 < f_2$ then the minimum lies between x_2 and x_3 and the new triplet is x_2, x_4 and x_3 .



Then

We want

$$x_3 - x_2 = \varphi(x_3 - x_1) \mapsto -x_2 = \varphi x_3 - \varphi x_1 - x_3$$

Therefore

$$x_2 = \varphi x_1 + (1 - \varphi)x_3$$

Thus, once we obtain x_1 , we get x_2 and x_3 .

- For this, we make the following assumption $[x_1, x_3] = [0, 1]$



Then

We want

$$x_3 - x_2 = \varphi(x_3 - x_1) \mapsto -x_2 = \varphi x_3 - \varphi x_1 - x_3$$

Therefore

$$x_2 = \varphi x_1 + (1 - \varphi)x_3$$

Thus, once we obtain x_1 , we get x_2 and x_3 .

- For this, we make the following assumption $[x_1, x_3] = [0, 1]$



Then

We want

$$x_3 - x_2 = \varphi(x_3 - x_1) \mapsto -x_2 = \varphi x_3 - \varphi x_1 - x_3$$

Therefore

$$x_2 = \varphi x_1 + (1 - \varphi)x_3$$

Thus, once we obtain φ , we get x_2 and x_4

- For this, we make the following assumption $[x_1, x_3] = [0, 1]$



Therefore

If we have $f_2 < f_4$

$$x_2 = 1 - \varphi$$

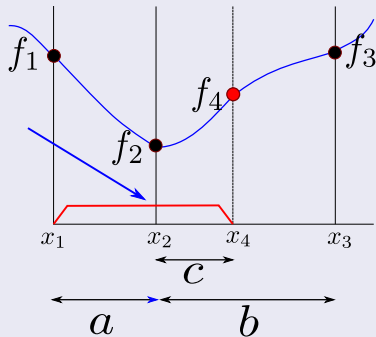
Then, if we have the new function evaluation at the left of x_2 ,

Therefore

If we have $f_2 < f_4$

$$x_2 = 1 - \varphi$$

Then, if we have the new function evaluation at the left of x_2



With a Little Algebra

Then, x_2 is between the the interval $[0, \varphi]$ and assume is a convex combination of such values

$$1 - \varphi = (1 - \varphi)0 + \varphi\varphi \mapsto \varphi^2 + \varphi - 1 = 0$$

With Solution

$$\varphi = \frac{-1 + \sqrt{5}}{2} = 0.6180$$



With a Little Algebra

Then, x_2 is between the the interval $[0, \varphi]$ and assume is a convex combination of such values

$$1 - \varphi = (1 - \varphi)0 + \varphi\varphi \mapsto \varphi^2 + \varphi - 1 = 0$$

With Solution

$$\varphi = \frac{-1 + \sqrt{5}}{2} = 0.6180$$



Finally, we have the algorithm

Golden Ratio

INPUT: $x_1, x_3, \tau, \varphi, f$

OUTPUT: $\frac{x_3 - x_1}{2}$

- 1 $x_2 = \varphi x_1 + (1 - \varphi)x_3$
- 2 $x_4 = (1 - \varphi)x_1 + \varphi x_3$
- 3 while $|x_3 - x_1| > \tau(|x_2| + |x_4|)$
- 4 if $f(x_2) < f(x_4)$:
- 5 $x_3 = x_4$
- 6 $x_4 = x_2$
- 7 $x_2 = \varphi x_1 + (1 - \varphi)x_3$
- 8 else
- 9 $x_1 = x_2$
- 10 $x_2 = x_4$
- 11 $x_4 = (1 - \varphi)x_1 + \varphi x_3$
- 12 return $\frac{x_3 - x_1}{2}$

Finally, we have the algorithm

Golden Ratio

INPUT: $x_1, x_3, \tau, \varphi, f$

OUTPUT: $\frac{x_3 - x_1}{2}$

- 1 $x_2 = \varphi x_1 + (1 - \varphi)x_3$
- 2 $x_4 = (1 - \varphi)x_1 + \varphi x_3$
- 3 while $|x_3 - x_1| > \tau(|x_2| + |x_4|)$
- 4 if $f(x_2) < f(x_4)$:
- 5 $x_3 = x_2$
- 6 $x_4 = x_4$
- 7 $x_2 = \varphi x_1 + (1 - \varphi)x_3$
- 8 else
- 9 $x_1 = x_2$
- 10 $x_2 = x_4$
- 11 $x_4 = (1 - \varphi)x_1 + \varphi x_3$
- 12 return $\frac{x_3 - x_1}{2}$

Finally, we have the algorithm

Golden Ratio

INPUT: $x_1, x_3, \tau, \varphi, f$

OUTPUT: $\frac{x_3 - x_1}{2}$

- 1 $x_2 = \varphi x_1 + (1 - \varphi)x_3$
- 2 $x_4 = (1 - \varphi)x_1 + \varphi x_3$
- 3 while $|x_3 - x_1| > \tau (|x_2| + |x_4|)$

 if $f(x_2) < f(x_4)$:

$x_3 = x_1$

$x_4 = x_2$

$x_2 = \varphi x_1 + (1 - \varphi)x_3$

 else

$x_1 = x_2$

$x_2 = x_4$

$x_4 = (1 - \varphi)x_1 + \varphi x_3$

return $\frac{x_3 - x_1}{2}$

Finally, we have the algorithm

Golden Ratio

INPUT: $x_1, x_3, \tau, \varphi, f$

OUTPUT: $\frac{x_3 - x_1}{2}$

- 1 $x_2 = \varphi x_1 + (1 - \varphi)x_3$
- 2 $x_4 = (1 - \varphi)x_1 + \varphi x_3$
- 3 while $|x_3 - x_1| > \tau (|x_2| + |x_4|)$
- 4 if $f(x_2) < f(x_4)$:
- 5 $x_3 = x_4$
- 6 $x_4 = x_2$
- 7 $x_2 = \varphi x_1 + (1 - \varphi)x_3$
- 8 else
- 9 $x_1 = x_2$
- 10 $x_2 = x_4$
- 11 $x_4 = (1 - \varphi)x_1 + \varphi x_3$
- 12 return $\frac{x_3 - x_1}{2}$

Finally, we have the algorithm

Golden Ratio

INPUT: $x_1, x_3, \tau, \varphi, f$

OUTPUT: $\frac{x_3 - x_1}{2}$

- 1 $x_2 = \varphi x_1 + (1 - \varphi)x_3$
- 2 $x_4 = (1 - \varphi)x_1 + \varphi x_3$
- 3 while $|x_3 - x_1| > \tau (|x_2| + |x_4|)$
- 4 if $f(x_2) < f(x_4)$:
- 5 $x_3 = x_4$
- 6 $x_4 = x_2$
- 7 $x_2 = \varphi x_1 + (1 - \varphi)x_3$
- 8 else
- 9 $x_1 = x_2$
- 10 $x_2 = x_4$
- 11 $x_4 = (1 - \varphi)x_1 + \varphi x_3$

return $\frac{x_3 - x_1}{2}$

Finally, we have the algorithm

Golden Ratio

INPUT: $x_1, x_3, \tau, \varphi, f$

OUTPUT: $\frac{x_3 - x_1}{2}$

- 1 $x_2 = \varphi x_1 + (1 - \varphi)x_3$
- 2 $x_4 = (1 - \varphi)x_1 + \varphi x_3$
- 3 while $|x_3 - x_1| > \tau (|x_2| + |x_4|)$
- 4 if $f(x_2) < f(x_4)$:
- 5 $x_3 = x_4$
- 6 $x_4 = x_2$
- 7 $x_2 = \varphi x_1 + (1 - \varphi)x_3$
- 8 else
- 9 $x_1 = x_2$
- 10 $x_2 = x_4$
- 11 $x_4 = (1 - \varphi)x_1 + \varphi x_3$
- 12 return $\frac{x_3 - x_1}{2}$

Iteratively

Repeat the procedure!!!

Until a error threshold is reached.



Iteratively

Repeat the procedure!!!

Until a error threshold is reached.

For more, please read the paper

“SEQUENTIAL MINIMAX SEARCH FOR A MAXIMUM” by J. Kiefer



Iteratively

Repeat the procedure!!!

Until a error threshold is reached.

For more, please read the paper

“SEQUENTIAL MINIMAX SEARCH FOR A MAXIMUM” by J. Kiefer



There are better versions

Take a look

The papers at the repository.

Outline

1 More in Regularization

- Introduction
- Smoothness of the Estimation
- The Error Estimate
- Choosing approximate inverses
- A Classic Example, Regularization as a Filter
- Another Example, The Landweber Iteration

2 Linear Regression using Gradient Descent

- Introduction
- What is the Gradient of the Equation?
- The Basic Algorithm
- How to obtain $\eta(k)$
 - Gold Section

3 The Gauss-Markov Theorem

- **Statement**
- Proof

4 Fisher Linear Discriminant

- History
- The Projection and The Rotation Idea
- Classifiers as Machines for dimensionality reduction
- Solution
 - Use the mean of each Class
 - Scatter measure
- The Cost Function
- A Transformation for simplification and defining the cost function
- Where is this used?
 - Applications
- Relation with Least Squared Error
 - What?

5 Exercises

- Some Stuff for you to try



The Gauss-Markov Theorem

Given the Linear Estimation Model

$$\mathbf{y} = \mathbf{X}\mathbf{w} + \boldsymbol{\epsilon}$$

Under the following assumptions

- 1. $E[\boldsymbol{\epsilon}|\mathbf{x}] = \mathbf{0}$ for all \mathbf{x} (Mean Independence).
- 2. $Var[\boldsymbol{\epsilon}] = E[\boldsymbol{\epsilon}\boldsymbol{\epsilon}^T|\mathbf{x}] = \sigma_c^2 \mathbf{I}_N$ (Homoskedasticity).

The Gauss-Markov Theorem states

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

is the Best Linear Unbiased Estimator (BLUE), if $\boldsymbol{\epsilon}$ satisfies 1. and 2.!!!



The Gauss-Markov Theorem

Given the Linear Estimation Model

$$\mathbf{y} = \mathbf{X}\mathbf{w} + \boldsymbol{\epsilon}$$

Under the following assumptions

- 1 $E[\boldsymbol{\epsilon}|\mathbf{x}] = \mathbf{0}$ for all \mathbf{x} (Mean Independence).
- 2 $Var[\boldsymbol{\epsilon}] = E[\boldsymbol{\epsilon}\boldsymbol{\epsilon}^T|\mathbf{x}] = \sigma_\epsilon^2 I_N$ (Homoskedasticity).

The Gauss-Markov Theorem states

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

is the Best Linear Unbiased Estimator (BLUE), if $\boldsymbol{\epsilon}$ satisfies 1. and 2.!!!



The Gauss-Markov Theorem

Given the Linear Estimation Model

$$\mathbf{y} = X\mathbf{w} + \boldsymbol{\epsilon}$$

Under the following assumptions

- 1 $E[\boldsymbol{\epsilon}|\mathbf{x}] = \mathbf{0}$ for all \mathbf{x} (Mean Independence).
- 2 $Var[\boldsymbol{\epsilon}] = E[\boldsymbol{\epsilon}\boldsymbol{\epsilon}^T|\mathbf{x}] = \sigma_\epsilon^2 I_N$ (Homoskedasticity).

The Gauss-Markov Theorem states

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

is the Best Linear Unbiased Estimator (BLUE), if $\boldsymbol{\epsilon}$ satisfies 1. and 2.!!!



Outline

1 More in Regularization

- Introduction
- Smoothness of the Estimation
- The Error Estimate
- Choosing approximate inverses
- A Classic Example, Regularization as a Filter
- Another Example, The Landweber Iteration

2 Linear Regression using Gradient Descent

- Introduction
- What is the Gradient of the Equation?
- The Basic Algorithm
- How to obtain $\eta(k)$
 - Gold Section

3 The Gauss-Markov Theorem

- Statement
- **Proof**

4 Fisher Linear Discriminant

- History
- The Projection and The Rotation Idea
- Classifiers as Machines for dimensionality reduction
- Solution
 - Use the mean of each Class
 - Scatter measure
- The Cost Function
- A Transformation for simplification and defining the cost function
- Where is this used?
 - Applications
- Relation with Least Squared Error
 - What?

5 Exercises

- Some Stuff for you to try



Proof

First and Fore most

- “An estimator is “best” in a class if it has smaller variance than others estimators in the same class.”

Also

- We are restricting our search for estimators to the class of linear, unbiased ones

Unbiased Estimator

Given a sequence of observations $x_1, x_2, \dots, x_N \sim P(X|\theta)$ then bias is the mean of the difference

$$b_d(\theta) = E[d(X) - h(\theta)]$$

with $d(X)$ is an estimator of the statistic $h(\theta)$.

Proof

First and Fore most

- “An estimator is “best” in a class if it has smaller variance than others estimators in the same class.”

Also

- We are restricting our search for estimators to the class of linear, unbiased ones

Unbiased Estimator

Given a sequence of observations $x_1, x_2, \dots, x_N \sim P(X|\theta)$ then bias is the mean of the difference

$$b_d(\theta) = E[d(X) - h(\theta)]$$

with $d(X)$ is an estimator of the statistic $h(\theta)$.

Proof

First and Fore most

- “An estimator is “best” in a class if it has smaller variance than others estimators in the same class.”

Also

- We are restricting our search for estimators to the class of linear, unbiased ones

Unbiased Estimator

Given a sequence of observations $x_1, x_2, \dots, x_N \sim P(X|\theta)$ then bias is the mean of the difference

$$b_d(\theta) = E[d(X) - h(\theta)]$$

with $d(X)$ is an estimator of the statistic $h(\theta)$.

Remark

We need to calculate estimators which have covariances

- The best estimator in a class of estimators is the one with **the “smallest” covariance matrix**

This

- We will look at such covariance matrix for the BLUE estimator.



Remark

We need to calculate estimators which have covariances

- The best estimator in a class of estimators is the one with **the “smallest” covariance matrix**

Thus

- We will look at such covariance matrix for the BLUE estimator.



Therefore, going back to our unbiased estimators

If $b_d(\theta) = 0$ for all values of the parameter

- Then, $d(X)$ is called an unbiased estimator.

Now, the data are the y , we are looking at estimators that are linear functions of y :

$$\tilde{w} = m + My$$

Here:

- \tilde{w} is a $k \times 1$ parameter vector
- m is a $k \times 1$ vector of constants,
- M is a $k \times N$ matrix of constants,
- The data vector y is $N \times 1$.

Therefore, going back to our unbiased estimators

If $b_d(\theta) = 0$ for all values of the parameter

- Then, $d(X)$ is called an unbiased estimator.

Now, the data are the \mathbf{y} , we are looking at estimators that are linear functions of \mathbf{y}

$$\tilde{\mathbf{w}} = \mathbf{m} + M\mathbf{y}$$

- $\tilde{\mathbf{w}}$ is a $k \times 1$ parameter vector
- \mathbf{m} is a $k \times 1$ vector of constants,
- M is a $k \times N$ matrix of constants,
- The data vector \mathbf{y} is $N \times 1$.

Therefore, going back to our unbiased estimators

If $b_d(\theta) = 0$ for all values of the parameter

- Then, $d(X)$ is called an unbiased estimator.

Now, the data are the \mathbf{y} , we are looking at estimators that are linear functions of \mathbf{y}

$$\tilde{\mathbf{w}} = \mathbf{m} + M\mathbf{y}$$

Here

- $\tilde{\mathbf{w}}$ is a $k \times 1$ parameter vector
- \mathbf{m} is a $k \times 1$ vector of constants,
- M is a $k \times N$ matrix of constants,
- The data vector \mathbf{y} is $N \times 1$.

Now

We are looking at unbiased estimators

$$E[\tilde{w}] = w$$



Now

We are looking at unbiased estimators

$$E[\tilde{w}] = w$$

if \tilde{w} is to be unbiased

$$\begin{aligned} E[\tilde{w}|X] &= m + ME[y|X] \\ &= m + ME[Xw + \epsilon|X] \\ &= m + MXw \end{aligned}$$



Now

We are looking at unbiased estimators

$$E[\tilde{w}] = w$$

if \tilde{w} is to be unbiased

$$\begin{aligned} E[\tilde{w}|X] &= m + ME[y|X] \\ &= m + ME[Xw + \epsilon|X] \\ &= m + MXw \end{aligned}$$



Cinvestav

Now

We are looking at unbiased estimators

$$E[\tilde{\mathbf{w}}] = \mathbf{w}$$

if $\tilde{\mathbf{w}}$ is to be unbiased

$$\begin{aligned} E[\tilde{\mathbf{w}}|X] &= \mathbf{m} + ME[\mathbf{y}|X] \\ &= \mathbf{m} + ME[X\mathbf{w} + \boldsymbol{\epsilon}|X] \\ &= \mathbf{m} + MX\mathbf{w} \end{aligned}$$



Now, we are forced

Given that we are looking for an unbiased estimator

$$\mathbf{m} = \mathbf{0} \text{ with } MX = I_k$$

For the least squared error

$$M = (X^T X)^{-1} X^T \implies MX = (X^T X)^{-1} X^T X = I_k$$

Looking for linear unbiased estimators requires to look for estimators as

$$\tilde{\mathbf{w}} = My \text{ with } MX = I_k$$



Now, we are forced

Given that we are looking for an unbiased estimator

$$m = 0 \text{ with } MX = I_k$$

For the least squared error

$$M = (X^T X)^{-1} X^T \implies MX = (X^T X)^{-1} X^T X = I_k$$

Looking for linear unbiased estimators requires to look for estimators

as

$$\tilde{w} = My \text{ with } MX = I_k$$



Now, we are forced

Given that we are looking for an unbiased estimator

$$\mathbf{m} = \mathbf{0} \text{ with } MX = I_k$$

For the least squared error

$$M = (X^T X)^{-1} X^T \iff MX = (X^T X)^{-1} X^T X = I_k$$

Looking for linear unbiased estimators requires to look for estimators as

$$\tilde{\mathbf{w}} = M\mathbf{y} \text{ with } MX = I_k$$



Therefore

We are looking at matrices as

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

where C is some $k \times n$ matrix.



Therefore

We are looking at matrices as

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

where C is some $k \times n$ matrix.

Now

$$\begin{aligned} MX &= \left[(X^T X)^{-1} X^T + C \right] X \\ &= I_k + CX = I_k \\ &\implies CX = 0 \end{aligned}$$



Therefore

We are looking at matrices as

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

where C is some $k \times n$ matrix.

Now

$$MX = \left[(X^T X)^{-1} X^T + C \right] X$$

$$= I_k + CX = I_k$$

$$\Rightarrow CX = 0$$



Therefore

We are looking at matrices as

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

where C is some $k \times n$ matrix.

Now

$$\begin{aligned} MX &= \left[(X^T X)^{-1} X^T + C \right] X \\ &= I_k + CX = I_k \end{aligned}$$

$$\Rightarrow CX = 0$$



Therefore

We are looking at matrices as

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

where C is some $k \times n$ matrix.

Now

$$\begin{aligned} MX &= \left[(X^T X)^{-1} X^T + C \right] X \\ &= I_k + CX = I_k \\ &\implies CX = 0 \end{aligned}$$



Therefore, we can compute the covariance matrix

For all alternative estimators \tilde{w}

$$\begin{aligned}\tilde{w} &= My \\ &= M[Xw + \epsilon] \\ &= w + M\epsilon\end{aligned}$$



Therefore, we can compute the covariance matrix

For all alternative estimators \tilde{w}

$$\begin{aligned}\tilde{w} &= My \\ &= M[Xw + \epsilon] \\ &= w + M\epsilon\end{aligned}$$

Therefore, the difference is $\tilde{w} - w = M\epsilon$

- And since the \tilde{w} is unbiased, $E[\tilde{w} - w|X] = 0$



Therefore, we can compute the covariance matrix

For all alternative estimators \tilde{w}

$$\begin{aligned}\tilde{w} &= My \\ &= M[Xw + \epsilon] \\ &= w + M\epsilon\end{aligned}$$

Therefore, the difference is $\tilde{w} - w = M\epsilon$

- And since the \tilde{w} is unbiased, $E[\tilde{w} - w|X] = 0$



Therefore, we can compute the covariance matrix

For all alternative estimators $\tilde{\mathbf{w}}$

$$\begin{aligned}\tilde{\mathbf{w}} &= M\mathbf{y} \\ &= M[X\mathbf{w} + \boldsymbol{\epsilon}] \\ &= \mathbf{w} + M\boldsymbol{\epsilon}\end{aligned}$$

Therefore, the difference is $\tilde{\mathbf{w}} - \mathbf{w} = M\boldsymbol{\epsilon}$

- And since the $\tilde{\mathbf{w}}$ is unbiased, $E[\tilde{\mathbf{w}} - \mathbf{w}|X] = 0$



We have

The Covariance Matrix

$$\begin{aligned} E \left[(\tilde{\mathbf{w}} - \mathbf{w}) (\tilde{\mathbf{w}} - \mathbf{w})^T | X \right] &= E \left[M \boldsymbol{\epsilon} (M \boldsymbol{\epsilon})^T | X \right] \\ &= E \left[M \boldsymbol{\epsilon} \boldsymbol{\epsilon}^T M^T | X \right] \\ &= M E \left[\boldsymbol{\epsilon} \boldsymbol{\epsilon}^T | X \right] M^T \\ &= M \sigma_{\epsilon}^2 I_N M^T \\ &= \sigma_{\epsilon}^2 M M^T \end{aligned}$$



We have

The Covariance Matrix

$$\begin{aligned} E \left[(\tilde{\mathbf{w}} - \mathbf{w}) (\tilde{\mathbf{w}} - \mathbf{w})^T | X \right] &= E \left[M \boldsymbol{\epsilon} (M \boldsymbol{\epsilon})^T | X \right] \\ &= E \left[M \boldsymbol{\epsilon} \boldsymbol{\epsilon}^T M^T | X \right] \\ &= M E \left[\boldsymbol{\epsilon} \boldsymbol{\epsilon}^T | X \right] M^T \\ &= M \sigma_c^2 I_N M^T \\ &= \sigma_c^2 M M^T \end{aligned}$$



We have

The Covariance Matrix

$$\begin{aligned} E \left[(\tilde{\mathbf{w}} - \mathbf{w}) (\tilde{\mathbf{w}} - \mathbf{w})^T | X \right] &= E \left[M \boldsymbol{\epsilon} (M \boldsymbol{\epsilon})^T | X \right] \\ &= E \left[M \boldsymbol{\epsilon} \boldsymbol{\epsilon}^T M^T | X \right] \\ &= M E \left[\boldsymbol{\epsilon} \boldsymbol{\epsilon}^T | X \right] M^T \\ &= M \sigma_c^2 I_N M^T \\ &= \sigma_c^2 M M^T \end{aligned}$$



We have

The Covariance Matrix

$$\begin{aligned} E \left[(\tilde{\mathbf{w}} - \mathbf{w}) (\tilde{\mathbf{w}} - \mathbf{w})^T | X \right] &= E \left[M \boldsymbol{\epsilon} (M \boldsymbol{\epsilon})^T | X \right] \\ &= E \left[M \boldsymbol{\epsilon} \boldsymbol{\epsilon}^T M^T | X \right] \\ &= M E \left[\boldsymbol{\epsilon} \boldsymbol{\epsilon}^T | X \right] M^T \\ &= M \sigma_{\epsilon}^2 I_N M^T \\ &= \sigma_{\epsilon}^2 M M^T \end{aligned}$$



We have

The Covariance Matrix

$$\begin{aligned} E \left[(\tilde{\mathbf{w}} - \mathbf{w}) (\tilde{\mathbf{w}} - \mathbf{w})^T | X \right] &= E \left[M \boldsymbol{\epsilon} (M \boldsymbol{\epsilon})^T | X \right] \\ &= E \left[M \boldsymbol{\epsilon} \boldsymbol{\epsilon}^T M^T | X \right] \\ &= M E \left[\boldsymbol{\epsilon} \boldsymbol{\epsilon}^T | X \right] M^T \\ &= M \sigma_{\epsilon}^2 I_N M^T \\ &= \sigma_{\epsilon}^2 M M^T \end{aligned}$$



Finally

Given that $CX = 0$

$$\begin{aligned}MM^T &= \left[(X^T X)^{-1} X^T + C \right] \left[(X^T X)^{-1} X^T + C \right]^T \\ &= (X^T X)^{-1} X^T X (X^T X)^{-1} + (X^T X)^{-1} X^T C \\ &\quad + C X (X^T X)^{-1} + C C^T \\ &= (X^T X)^{-1} + C C^T\end{aligned}$$



Finally

Given that $CX = 0$

$$\begin{aligned}MM^T &= \left[(X^T X)^{-1} X^T + C \right] \left[(X^T X)^{-1} X^T + C \right]^T \\ &= (X^T X)^{-1} X^T X (X^T X)^{-1} + (X^T X)^{-1} X^T C \\ &\quad + CX (X^T X)^{-1} + CC^T \\ &= (X^T X)^{-1} + CC^T\end{aligned}$$

Now the matrix CC^T is a $k \times k$ "cross products" matrix

- By construction is positive semi-definite



Finally

Given that $CX = 0$

$$\begin{aligned}MM^T &= \left[(X^T X)^{-1} X^T + C \right] \left[(X^T X)^{-1} X^T + C \right]^T \\ &= (X^T X)^{-1} X^T X (X^T X)^{-1} + (X^T X)^{-1} X^T C \\ &\quad + CX (X^T X)^{-1} + CC^T \\ &= (X^T X)^{-1} + CC^T\end{aligned}$$

Now the matrix CC^T is a $k \times k$ "cross products" matrix

- By construction is positive semi-definite



Finally

Given that $CX = 0$

$$\begin{aligned}MM^T &= \left[(X^T X)^{-1} X^T + C \right] \left[(X^T X)^{-1} X^T + C \right]^T \\ &= (X^T X)^{-1} X^T X (X^T X)^{-1} + (X^T X)^{-1} X^T C \\ &\quad + CX (X^T X)^{-1} + CC^T \\ &= (X^T X)^{-1} + CC^T\end{aligned}$$

Now the matrix CC^T is a $k \times k$ “cross products” matrix

- By construction is positive semi-definite



Thus

Given

- The best estimator in a class of estimators is the one with the “smallest” covariance matrix

Where by “small”

- The covariance matrix associated with any other estimator in the class minus the covariance matrix of the best estimator is a positive definite matrix



Thus

Given

- The best estimator in a class of estimators is the one with the “smallest” covariance matrix

Where by “small”

- The covariance matrix associated with any other estimator in the class minus the covariance matrix of the best estimator is a **positive definite matrix**



Formally

The following difference is positive definite

$$MM^T + CC^T - Cov_{best}$$



Then

Since $MM^T + CC^T - Cov_{best}$ is minimized when we set the matrix C equal to the 0 matrix

- i.e. $M = (X^T X)^{-1} X$
 - ▶ The best estimator in the class \hat{w} .

Any other estimator \tilde{w} in this class

- It has strictly “larger” covariance matrix

Therefore the Least Square Error estimator \hat{w}

- It is BLUE under the two conditions of mean independence and homoskedastic!!!



Then

Since $MM^T + CC^T - Cov_{best}$ is minimized when we set the matrix C equal to the 0 matrix

- i.e. $M = (X^T X)^{-1} X$
 - ▶ The best estimator in the class \hat{w} .

Any other estimator M in this class

- It has strictly “larger” covariance matrix

Therefore, the Least Squares Error estimator \hat{w}

- It is BLUE under the two conditions of mean independence and homoskedastic!!!



Then

Since $MM^T + CC^T - Cov_{best}$ is minimized when we set the matrix C equal to the 0 matrix

- i.e. $M = (X^T X)^{-1} X$
 - ▶ The best estimator in the class \hat{w} .

Any other estimator M in this class

- It has strictly “larger” covariance matrix

Therefore the Least Square Error estimator \hat{w}

- It is BLUE under the two conditions of mean independence and homoskedastic!!!



Outline

1 More in Regularization

- Introduction
- Smoothness of the Estimation
- The Error Estimate
- Choosing approximate inverses
- A Classic Example, Regularization as a Filter
- Another Example, The Landweber Iteration

2 Linear Regression using Gradient Descent

- Introduction
- What is the Gradient of the Equation?
- The Basic Algorithm
- How to obtain $\eta(k)$
 - Gold Section

3 The Gauss-Markov Theorem

- Statement
- Proof

4 Fisher Linear Discriminant

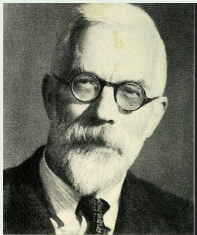
- **History**
- The Projection and The Rotation Idea
- Classifiers as Machines for dimensionality reduction
- Solution
 - Use the mean of each Class
 - Scatter measure
- The Cost Function
- A Transformation for simplification and defining the cost function
- Where is this used?
 - Applications
- Relation with Least Squared Error
 - What?

5 Exercises

- Some Stuff for you to try

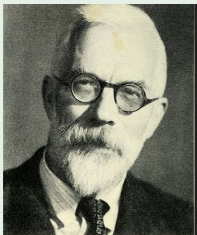
Invented Originally by

Sir Ronald Fisher



Invented Originally by

Sir Ronald Fisher

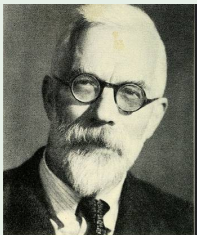


Anders Hald called him

"A genius who almost single-handedly created the foundations for modern statistical science."

Invented Originally by

Sir Ronald Fisher



Anders Hald called him

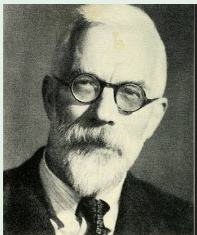
"A genius who almost single-handedly created the foundations for modern statistical science."

The Darkest Side

- In 1910 he joined the Eugenics Society at Cambridge, whose members included John Maynard Keynes, R. C. Punnett, and Horace Darwin.
- He opposed UNESCO's *The Race Question*, believing that evidence and everyday experience showed that human groups differ profoundly.

Invented Originally by

Sir Ronald Fisher



Anders Hald called him

"A genius who almost single-handedly created the foundations for modern statistical science."

The Darkest Side

- In 1910 he joined the Eugenics Society at Cambridge, whose members included John Maynard Keynes, R. C. Punnett, and Horace Darwin.
- He opposed UNESCO's The Race Question, believing that evidence and everyday experience showed that human groups differ profoundly.

Outline

1 More in Regularization

- Introduction
- Smoothness of the Estimation
- The Error Estimate
- Choosing approximate inverses
- A Classic Example, Regularization as a Filter
- Another Example, The Landweber Iteration

2 Linear Regression using Gradient Descent

- Introduction
- What is the Gradient of the Equation?
- The Basic Algorithm
- How to obtain $\eta(k)$
 - Gold Section

3 The Gauss-Markov Theorem

- Statement
- Proof

4 Fisher Linear Discriminant

- History
- **The Projection and The Rotation Idea**
- Classifiers as Machines for dimensionality reduction
- Solution
 - Use the mean of each Class
 - Scatter measure
- The Cost Function
- A Transformation for simplification and defining the cost function
- Where is this used?
 - Applications
- Relation with Least Squared Error
 - What?

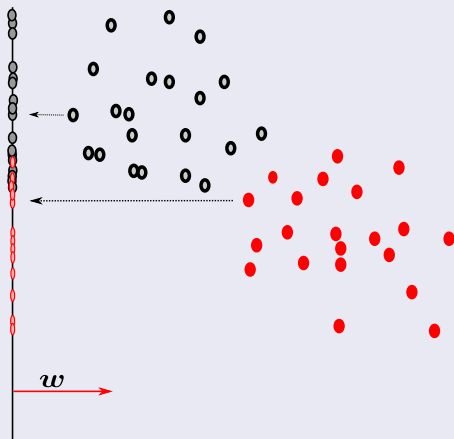
5 Exercises

- Some Stuff for you to try



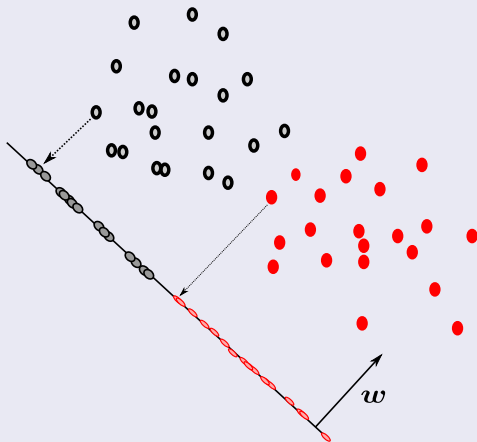
Intuition

Something Notable - Projecting into a Line



A Better Line

Something Notable - Projecting into a Line



Rotation

Projecting

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

Something Notable

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

Principal Component Analysis (PCA)

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



Rotation

Projecting

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

Something Notable

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

Principal Component Analysis (PCA)

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



Rotation

Projecting

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

Something Notable

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

Fisher Linear Discriminant (FLD)

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



Outline

1 More in Regularization

- Introduction
- Smoothness of the Estimation
- The Error Estimate
- Choosing approximate inverses
- A Classic Example, Regularization as a Filter
- Another Example, The Landweber Iteration

2 Linear Regression using Gradient Descent

- Introduction
- What is the Gradient of the Equation?
- The Basic Algorithm
- How to obtain $\eta(k)$
 - Gold Section

3 The Gauss-Markov Theorem

- Statement
- Proof

4 Fisher Linear Discriminant

- History
- The Projection and The Rotation Idea
- **Classifiers as Machines for dimensionality reduction**
- Solution
 - Use the mean of each Class
 - Scatter measure
- The Cost Function
- A Transformation for simplification and defining the cost function
- Where is this used?
 - Applications
- Relation with Least Squared Error
 - What?

5 Exercises

- Some Stuff for you to try



This is actually coming from...

Classifier as

A machine for dimensionality reduction.



Cinvestav

This is actually coming from...

Classifier as

A machine for dimensionality reduction.

Initial Setup

We have:

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



This is actually coming from...

Classifier as

A machine for dimensionality reduction.

Initial Setup

We have:

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

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

$$y_i = w^T x_i \quad (9)$$



This is actually coming from...

Classifier as

A machine for dimensionality reduction.

Initial Setup

We have:

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

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

$$y_i = w^T x_i \quad (9)$$



This is actually coming from...

Classifier as

A machine for dimensionality reduction.

Initial Setup

We have:

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

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

$$y_i = \mathbf{w}^T \mathbf{x}_i \quad (9)$$



Outline

1 More in Regularization

- Introduction
- Smoothness of the Estimation
- The Error Estimate
- Choosing approximate inverses
- A Classic Example, Regularization as a Filter
- Another Example, The Landweber Iteration

2 Linear Regression using Gradient Descent

- Introduction
- What is the Gradient of the Equation?
- The Basic Algorithm
- How to obtain $\eta(k)$
 - Gold Section

3 The Gauss-Markov Theorem

- Statement
- Proof

4 Fisher Linear Discriminant

- History
- The Projection and The Rotation Idea
- Classifiers as Machines for dimensionality reduction
- **Solution**
 - Use the mean of each Class
 - Scatter measure
- The Cost Function
- A Transformation for simplification and defining the cost function
- Where is this used?
 - Applications
- Relation with Least Squared Error
 - What?

5 Exercises

- Some Stuff for you to try



Outline

1 More in Regularization

- Introduction
- Smoothness of the Estimation
- The Error Estimate
- Choosing approximate inverses
- A Classic Example, Regularization as a Filter
- Another Example, The Landweber Iteration

2 Linear Regression using Gradient Descent

- Introduction
- What is the Gradient of the Equation?
- The Basic Algorithm
- How to obtain $\eta(k)$
 - Gold Section

3 The Gauss-Markov Theorem

- Statement
- Proof

4 Fisher Linear Discriminant

- History
- The Projection and The Rotation Idea
- Classifiers as Machines for dimensionality reduction
- **Solution**
 - **Use the mean of each Class**
 - Scatter measure
- The Cost Function
- A Transformation for simplification and defining the cost function
- Where is this used?
 - Applications
- Relation with Least Squared Error
 - What?

5 Exercises

- Some Stuff for you to try



Use the mean of each Class

Then

Select w such that class separation is maximized

Use the mean of each Class

Then

Select w such that class separation is maximized

We then define the mean sample for each class

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

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

Use the mean of each Class

Then

Select w such that class separation is maximized

We then define the mean sample for each class

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

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

Goal: This is a linear measure of distance

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

$$m_1 - m_2 = w^T (m_1 - m_2) \quad (10)$$

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

Use the mean of each Class

Then

Select w such that class separation is maximized

We then define the mean sample for each class

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

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

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

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

$$m_1 - m_2 = w^T (m_1 - m_2) \quad (10)$$

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

Use the mean of each Class

Then

Select w such that class separation is maximized

We then define the mean sample for each class

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

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

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

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

$$m_1 - m_2 = w^T (m_1 - m_2) \quad (10)$$

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

However

We could simply seek

$$\begin{aligned} \max_w \quad & \mathbf{w}^T (\mathbf{m}_1 - \mathbf{m}_2) \\ \text{s.t.} \quad & \sqrt{\mathbf{w}^T \mathbf{w}} = 1 \end{aligned}$$

After all

We do not care about the magnitude of w .



However

We could simply seek

$$\begin{aligned} \max_w \quad & \mathbf{w}^T (\mathbf{m}_1 - \mathbf{m}_2) \\ \text{s.t.} \quad & \sqrt{\mathbf{w}^T \mathbf{w}} = 1 \end{aligned}$$

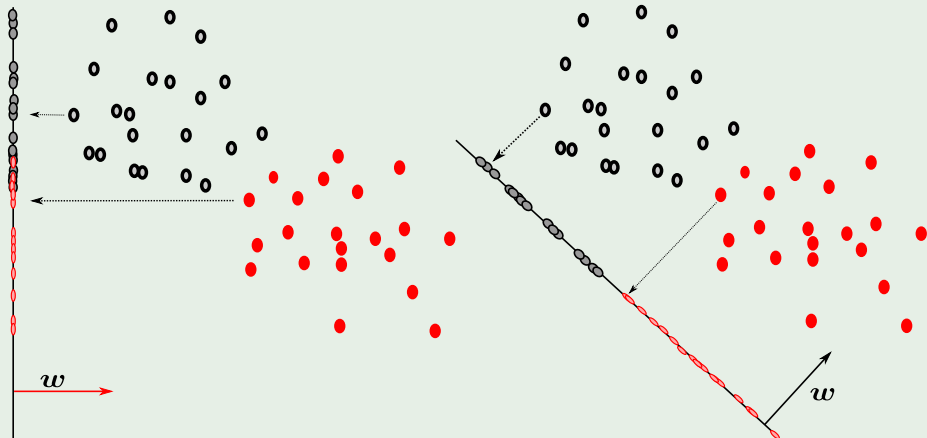
After all

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



Example

Here, we have the problem... The Scattering!!!



Outline

1 More in Regularization

- Introduction
- Smoothness of the Estimation
- The Error Estimate
- Choosing approximate inverses
- A Classic Example, Regularization as a Filter
- Another Example, The Landweber Iteration

2 Linear Regression using Gradient Descent

- Introduction
- What is the Gradient of the Equation?
- The Basic Algorithm
- How to obtain $\eta(k)$
 - Gold Section

3 The Gauss-Markov Theorem

- Statement
- Proof

4 Fisher Linear Discriminant

- History
- The Projection and The Rotation Idea
- Classifiers as Machines for dimensionality reduction
- **Solution**
 - Use the mean of each Class
 - **Scatter measure**
- The Cost Function
- A Transformation for simplification and defining the cost function
- Where is this used?
 - Applications
- Relation with Least Squared Error
 - What?

5 Exercises

- Some Stuff for you to try



Fixing the Problem

To obtain good separation of the projected data

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

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

$$s_k^2 = \sum_{x_i \in C_k} (w^T x_i - m_k)^2 = \sum_{y_i = w^T x_i \in C_k} (y_i - m_k)^2 \quad (11)$$

We define then within-class variance for the whole data

$$s_1^2 + s_2^2 \quad (12)$$



Fixing the Problem

To obtain good separation of the projected data

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

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

$$s_k^2 = \sum_{\mathbf{x}_i \in C_k} (\mathbf{w}^T \mathbf{x}_i - m_k)^2 = \sum_{y_i = \mathbf{w}^T \mathbf{x}_i \in C_k} (y_i - m_k)^2 \quad (11)$$

We define then within-class variance for the whole data

$$s_1^2 + s_2^2 \quad (12)$$



Fixing the Problem

To obtain good separation of the projected data

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

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

$$s_k^2 = \sum_{\mathbf{x}_i \in C_k} (\mathbf{w}^T \mathbf{x}_i - m_k)^2 = \sum_{y_i = \mathbf{w}^T \mathbf{x}_i \in C_k} (y_i - m_k)^2 \quad (11)$$

We define then within-class variance for the whole data

$$s_1^2 + s_2^2 \quad (12)$$



Outline

1 More in Regularization

- Introduction
- Smoothness of the Estimation
- The Error Estimate
- Choosing approximate inverses
- A Classic Example, Regularization as a Filter
- Another Example, The Landweber Iteration

2 Linear Regression using Gradient Descent

- Introduction
- What is the Gradient of the Equation?
- The Basic Algorithm
- How to obtain $\eta(k)$
 - Gold Section

3 The Gauss-Markov Theorem

- Statement
- Proof

4 Fisher Linear Discriminant

- History
- The Projection and The Rotation Idea
- Classifiers as Machines for dimensionality reduction
- Solution
 - Use the mean of each Class
 - Scatter measure
- **The Cost Function**
 - A Transformation for simplification and defining the cost function
 - Where is this used?
 - Applications
 - Relation with Least Squared Error
 - What?

5 Exercises

- Some Stuff for you to try



Finally, a Cost Function

The between-class variance

$$(m_1 - m_2)^2 \quad (13)$$

The Fisher Criterion (A Ratio)

$$\frac{\text{between-class variance}}{\text{within-class variance}} \quad (14)$$

Finally

$$J(w) = \frac{(m_1 - m_2)^2}{s_1^2 + s_2^2} \quad (15)$$



Finally, a Cost Function

The between-class variance

$$(m_1 - m_2)^2 \quad (13)$$

The Fisher criterion (A Ratio)

$$\frac{\text{between-class variance}}{\text{within-class variance}} \quad (14)$$

Finally

$$J(w) = \frac{(m_1 - m_2)^2}{s_1^2 + s_2^2} \quad (15)$$



Finally, a Cost Function

The between-class variance

$$(m_1 - m_2)^2 \quad (13)$$

The Fisher criterion (A Ratio)

$$\frac{\text{between-class variance}}{\text{within-class variance}} \quad (14)$$

Finally

$$J(\mathbf{w}) = \frac{(m_1 - m_2)^2}{s_1^2 + s_2^2} \quad (15)$$



Outline

1 More in Regularization

- Introduction
- Smoothness of the Estimation
- The Error Estimate
- Choosing approximate inverses
- A Classic Example, Regularization as a Filter
- Another Example, The Landweber Iteration

2 Linear Regression using Gradient Descent

- Introduction
- What is the Gradient of the Equation?
- The Basic Algorithm
- How to obtain $\eta(k)$
 - Gold Section

3 The Gauss-Markov Theorem

- Statement
- Proof

4 Fisher Linear Discriminant

- History
- The Projection and The Rotation Idea
- Classifiers as Machines for dimensionality reduction
- Solution
 - Use the mean of each Class
 - Scatter measure
- The Cost Function
- **A Transformation for simplification and defining the cost function**
- Where is this used?
 - Applications
- Relation with Least Squared Error
 - What?

5 Exercises

- Some Stuff for you to try



We use a transformation to simplify our life

First

$$J(\mathbf{w}) = \frac{(\mathbf{w}^T \mathbf{m}_1 - \mathbf{w}^T \mathbf{m}_2)^2}{\sum_{y_i = \mathbf{w}^T \mathbf{x}_i \in C_1} (y_i - m_1)^2 + \sum_{y_i = \mathbf{w}^T \mathbf{x}_i \in C_2} (y_i - m_2)^2} \quad (16)$$

Second

$$\frac{(\mathbf{w}^T \mathbf{m}_1 - \mathbf{w}^T \mathbf{m}_2) (\mathbf{w}^T \mathbf{m}_1 - \mathbf{w}^T \mathbf{m}_2)^T}{\sum_{y_i = \mathbf{w}^T \mathbf{x}_i \in C_1} (\mathbf{w}^T \mathbf{x}_i - m_1) (\mathbf{w}^T \mathbf{x}_i - m_1)^T + \sum_{y_i = \mathbf{w}^T \mathbf{x}_i \in C_2} (\mathbf{w}^T \mathbf{x}_i - m_2) (\mathbf{w}^T \mathbf{x}_i - m_2)^T} \quad (17)$$

Third

$$\frac{\mathbf{w}^T (\mathbf{m}_1 - \mathbf{m}_2) (\mathbf{w}^T (\mathbf{m}_1 - \mathbf{m}_2))^T}{\sum_{y_i = \mathbf{w}^T \mathbf{x}_i \in C_1} \mathbf{w}^T (\mathbf{x}_i - m_1) (\mathbf{w}^T (\mathbf{x}_i - m_1))^T + \sum_{y_i = \mathbf{w}^T \mathbf{x}_i \in C_2} \mathbf{w}^T (\mathbf{x}_i - m_2) (\mathbf{w}^T (\mathbf{x}_i - m_2))^T} \quad (18)$$

We use a transformation to simplify our life

First

$$J(\mathbf{w}) = \frac{(\mathbf{w}^T \mathbf{m}_1 - \mathbf{w}^T \mathbf{m}_2)^2}{\sum_{y_i = \mathbf{w}^T \mathbf{x}_i \in C_1} (y_i - m_1)^2 + \sum_{y_i = \mathbf{w}^T \mathbf{x}_i \in C_2} (y_i - m_2)^2} \quad (16)$$

Second

$$\frac{(\mathbf{w}^T \mathbf{m}_1 - \mathbf{w}^T \mathbf{m}_2) (\mathbf{w}^T \mathbf{m}_1 - \mathbf{w}^T \mathbf{m}_2)^T}{\sum_{y_i = \mathbf{w}^T \mathbf{x}_i \in C_1} (\mathbf{w}^T \mathbf{x}_i - m_1) (\mathbf{w}^T \mathbf{x}_i - m_1)^T + \sum_{y_i = \mathbf{w}^T \mathbf{x}_i \in C_2} (\mathbf{w}^T \mathbf{x}_i - m_2) (\mathbf{w}^T \mathbf{x}_i - m_2)^T} \quad (17)$$

Third

$$\frac{\mathbf{w}^T (\mathbf{m}_1 - \mathbf{m}_2) (\mathbf{w}^T (\mathbf{m}_1 - \mathbf{m}_2))^T}{\sum_{y_i = \mathbf{w}^T \mathbf{x}_i \in C_1} \mathbf{w}^T (\mathbf{x}_i - m_1) (\mathbf{w}^T (\mathbf{x}_i - m_1))^T + \sum_{y_i = \mathbf{w}^T \mathbf{x}_i \in C_2} \mathbf{w}^T (\mathbf{x}_i - m_2) (\mathbf{w}^T (\mathbf{x}_i - m_2))^T} \quad (18)$$

We use a transformation to simplify our life

First

$$J(\mathbf{w}) = \frac{(\mathbf{w}^T \mathbf{m}_1 - \mathbf{w}^T \mathbf{m}_2)^2}{\sum_{y_i = \mathbf{w}^T \mathbf{x}_i \in C_1} (y_i - m_1)^2 + \sum_{y_i = \mathbf{w}^T \mathbf{x}_i \in C_2} (y_i - m_2)^2} \quad (16)$$

Second

$$\frac{(\mathbf{w}^T \mathbf{m}_1 - \mathbf{w}^T \mathbf{m}_2) (\mathbf{w}^T \mathbf{m}_1 - \mathbf{w}^T \mathbf{m}_2)^T}{\sum_{y_i = \mathbf{w}^T \mathbf{x}_i \in C_1} (\mathbf{w}^T \mathbf{x}_i - m_1) (\mathbf{w}^T \mathbf{x}_i - m_1)^T + \sum_{y_i = \mathbf{w}^T \mathbf{x}_i \in C_2} (\mathbf{w}^T \mathbf{x}_i - m_2) (\mathbf{w}^T \mathbf{x}_i - m_2)^T} \quad (17)$$

Third

$$\frac{\mathbf{w}^T (\mathbf{m}_1 - \mathbf{m}_2) (\mathbf{w}^T (\mathbf{m}_1 - \mathbf{m}_2))^T}{\sum_{y_i = \mathbf{w}^T \mathbf{x}_i \in C_1} \mathbf{w}^T (\mathbf{x}_i - \mathbf{m}_1) (\mathbf{w}^T (\mathbf{x}_i - \mathbf{m}_1))^T + \sum_{y_i = \mathbf{w}^T \mathbf{x}_i \in C_2} \mathbf{w}^T (\mathbf{x}_i - \mathbf{m}_2) (\mathbf{w}^T (\mathbf{x}_i - \mathbf{m}_2))^T} \quad (18)$$

Transformation

Fourth

$$\frac{w^T (m_1 - m_2) (m_1 - m_2)^T w}{\sum_{y_i = w^T x_i \in C_1} w^T (x_i - m_1) (x_i - m_1)^T w + \sum_{y_i = w^T x_i \in C_2} w^T (x_i - m_2) (x_i - m_2)^T w} \quad (19)$$

Fifth

$$\frac{w^T (m_1 - m_2) (m_1 - m_2)^T w}{w^T \left[\sum_{y_i = w^T x_i \in C_1} (x_i - m_1) (x_i - m_1)^T + \sum_{y_i = w^T x_i \in C_2} (x_i - m_2) (x_i - m_2)^T \right] w} \quad (20)$$

Now Rename

$$J(w) = \frac{w^T S_B w}{w^T S_w w} \quad (21)$$



Transformation

Fourth

$$\frac{\mathbf{w}^T (\mathbf{m}_1 - \mathbf{m}_2) (\mathbf{m}_1 - \mathbf{m}_2)^T \mathbf{w}}{\sum_{y_i = \mathbf{w}^T \mathbf{x}_i \in C_1} \mathbf{w}^T (\mathbf{x}_i - \mathbf{m}_1) (\mathbf{x}_i - \mathbf{m}_1)^T \mathbf{w} + \sum_{y_i = \mathbf{w}^T \mathbf{x}_i \in C_2} \mathbf{w}^T (\mathbf{x}_i - \mathbf{m}_2) (\mathbf{x}_i - \mathbf{m}_2)^T \mathbf{w}} \quad (19)$$

Fifth

$$\frac{\mathbf{w}^T (\mathbf{m}_1 - \mathbf{m}_2) (\mathbf{m}_1 - \mathbf{m}_2)^T \mathbf{w}}{\mathbf{w}^T \left[\sum_{y_i = \mathbf{w}^T \mathbf{x}_i \in C_1} (\mathbf{x}_i - \mathbf{m}_1) (\mathbf{x}_i - \mathbf{m}_1)^T + \sum_{y_i = \mathbf{w}^T \mathbf{x}_i \in C_2} (\mathbf{x}_i - \mathbf{m}_2) (\mathbf{x}_i - \mathbf{m}_2)^T \right] \mathbf{w}} \quad (20)$$

Now Remain

$$J(\mathbf{w}) = \frac{\mathbf{w}^T \mathcal{S}_B \mathbf{w}}{\mathbf{w}^T \mathcal{S}_W \mathbf{w}} \quad (21)$$



Transformation

Fourth

$$\frac{\mathbf{w}^T (\mathbf{m}_1 - \mathbf{m}_2) (\mathbf{m}_1 - \mathbf{m}_2)^T \mathbf{w}}{\sum_{y_i = \mathbf{w}^T \mathbf{x}_i \in C_1} \mathbf{w}^T (\mathbf{x}_i - \mathbf{m}_1) (\mathbf{x}_i - \mathbf{m}_1)^T \mathbf{w} + \sum_{y_i = \mathbf{w}^T \mathbf{x}_i \in C_2} \mathbf{w}^T (\mathbf{x}_i - \mathbf{m}_2) (\mathbf{x}_i - \mathbf{m}_2)^T \mathbf{w}} \quad (19)$$

Fifth

$$\frac{\mathbf{w}^T (\mathbf{m}_1 - \mathbf{m}_2) (\mathbf{m}_1 - \mathbf{m}_2)^T \mathbf{w}}{\mathbf{w}^T \left[\sum_{y_i = \mathbf{w}^T \mathbf{x}_i \in C_1} (\mathbf{x}_i - \mathbf{m}_1) (\mathbf{x}_i - \mathbf{m}_1)^T + \sum_{y_i = \mathbf{w}^T \mathbf{x}_i \in C_2} (\mathbf{x}_i - \mathbf{m}_2) (\mathbf{x}_i - \mathbf{m}_2)^T \right] \mathbf{w}} \quad (20)$$

Now Rename

$$J(\mathbf{w}) = \frac{\mathbf{w}^T \mathbf{S}_B \mathbf{w}}{\mathbf{w}^T \mathbf{S}_w \mathbf{w}} \quad (21)$$



Derive with respect to \mathbf{w}

Thus

$$\frac{dJ(\mathbf{w})}{d\mathbf{w}} = \frac{d(\mathbf{w}^T \mathbf{S}_B \mathbf{w}) (\mathbf{w}^T \mathbf{S}_w \mathbf{w})^{-1}}{d\mathbf{w}} = 0 \quad (22)$$

Then

$$\frac{dJ(\mathbf{w})}{d\mathbf{w}} = (\mathbf{S}_B \mathbf{w} + \mathbf{S}_B^T \mathbf{w}) (\mathbf{w}^T \mathbf{S}_w \mathbf{w})^{-1} - (\mathbf{w}^T \mathbf{S}_B \mathbf{w}) (\mathbf{w}^T \mathbf{S}_w \mathbf{w})^{-2} (\mathbf{S}_w \mathbf{w} + \mathbf{S}_w^T \mathbf{w}) = 0 \quad (23)$$

Now, because the symmetry in \mathbf{S}_B and \mathbf{S}_w

$$\frac{dJ(\mathbf{w})}{d\mathbf{w}} = \frac{\mathbf{S}_B \mathbf{w}}{(\mathbf{w}^T \mathbf{S}_w \mathbf{w})} - \frac{\mathbf{w}^T \mathbf{S}_B \mathbf{w} \mathbf{S}_w \mathbf{w}}{(\mathbf{w}^T \mathbf{S}_w \mathbf{w})^2} = 0 \quad (24)$$



Derive with respect to w

Thus

$$\frac{dJ(w)}{dw} = \frac{d(w^T S_B w) (w^T S_w w)^{-1}}{dw} = 0 \quad (22)$$

Then

$$\frac{dJ(w)}{dw} = (S_B w + S_B^T w) (w^T S_w w)^{-1} - (w^T S_B w) (w^T S_w w)^{-2} (S_w w + S_w^T w) = 0 \quad (23)$$

Now, because the symmetry in S_B and S_w

$$\frac{dJ(w)}{dw} = \frac{S_B w}{(w^T S_w w)} - \frac{w^T S_B w S_w w}{(w^T S_w w)^2} = 0 \quad (24)$$

Derive with respect to w

Thus

$$\frac{dJ(w)}{dw} = \frac{d(w^T S_B w) (w^T S_w w)^{-1}}{dw} = 0 \quad (22)$$

Then

$$\frac{dJ(w)}{dw} = (S_B w + S_B^T w) (w^T S_w w)^{-1} - (w^T S_B w) (w^T S_w w)^{-2} (S_w w + S_w^T w) = 0 \quad (23)$$

Now because the symmetry in S_B and S_w

$$\frac{dJ(w)}{dw} = \frac{S_B w}{(w^T S_w w)} - \frac{w^T S_B w S_w w}{(w^T S_w w)^2} = 0 \quad (24)$$



Derive with respect to w

Thus

$$\frac{dJ(w)}{dw} = \frac{S_B w}{(w^T S_w w)} - \frac{w^T S_B w S_w w}{(w^T S_w w)^2} = 0 \quad (25)$$

Then

$$(w^T S_w w) S_B w = (w^T S_B w) S_w w \quad (26)$$



Derive with respect to w

Thus

$$\frac{dJ(w)}{dw} = \frac{S_B w}{(w^T S_w w)} - \frac{w^T S_B w S_w w}{(w^T S_w w)^2} = 0 \quad (25)$$

Then

$$(w^T S_w w) S_B w = (w^T S_B w) S_w w \quad (26)$$



Now, Several Tricks!!!

First

$$S_B w = (m_1 - m_2) (m_1 - m_2)^T w = \alpha (m_1 - m_2) \quad (27)$$

Where $\alpha = (m_1 - m_2)^T w$ is a simple constant

It means that $S_B w$ is always in the direction $m_1 - m_2$!!!

In addition

$w^T S_w w$ and $w^T S_B w$ are constants



Now, Several Tricks!!!

First

$$\mathbf{S}_B \mathbf{w} = (\mathbf{m}_1 - \mathbf{m}_2) (\mathbf{m}_1 - \mathbf{m}_2)^T \mathbf{w} = \alpha (\mathbf{m}_1 - \mathbf{m}_2) \quad (27)$$

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

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

In addition

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



Now, Several Tricks!!!

First

$$\mathbf{S}_B \mathbf{w} = (\mathbf{m}_1 - \mathbf{m}_2) (\mathbf{m}_1 - \mathbf{m}_2)^T \mathbf{w} = \alpha (\mathbf{m}_1 - \mathbf{m}_2) \quad (27)$$

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

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

In addition

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



Now, Several Tricks!!!

Finally, we only need the direction

$$\mathbf{S}_w \mathbf{w} \propto (\mathbf{m}_1 - \mathbf{m}_2) \Rightarrow \mathbf{w} \propto \mathbf{S}_w^{-1} (\mathbf{m}_1 - \mathbf{m}_2) \quad (28)$$



Now, Several Tricks!!!

Finally, we only need the direction

$$\mathbf{S}_w \mathbf{w} \propto (\mathbf{m}_1 - \mathbf{m}_2) \Rightarrow \mathbf{w} \propto \mathbf{S}_w^{-1} (\mathbf{m}_1 - \mathbf{m}_2) \quad (28)$$

Once the data is transformed into y_i

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



Now, Several Tricks!!!

Finally, we only need the direction

$$\mathbf{S}_w \mathbf{w} \propto (\mathbf{m}_1 - \mathbf{m}_2) \Rightarrow \mathbf{w} \propto \mathbf{S}_w^{-1} (\mathbf{m}_1 - \mathbf{m}_2) \quad (28)$$

Once the data is transformed into y_i

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



Outline

1 More in Regularization

- Introduction
- Smoothness of the Estimation
- The Error Estimate
- Choosing approximate inverses
- A Classic Example, Regularization as a Filter
- Another Example, The Landweber Iteration

2 Linear Regression using Gradient Descent

- Introduction
- What is the Gradient of the Equation?
- The Basic Algorithm
- How to obtain $\eta(k)$
 - Gold Section

3 The Gauss-Markov Theorem

- Statement
- Proof

4 Fisher Linear Discriminant

- History
- The Projection and The Rotation Idea
- Classifiers as Machines for dimensionality reduction
- Solution
 - Use the mean of each Class
 - Scatter measure
- The Cost Function
- A Transformation for simplification and defining the cost function
- **Where is this used?**
 - Applications
- Relation with Least Squared Error
 - What?

5 Exercises

- Some Stuff for you to try



Outline

1 More in Regularization

- Introduction
- Smoothness of the Estimation
- The Error Estimate
- Choosing approximate inverses
- A Classic Example, Regularization as a Filter
- Another Example, The Landweber Iteration

2 Linear Regression using Gradient Descent

- Introduction
- What is the Gradient of the Equation?
- The Basic Algorithm
- How to obtain $\eta(k)$
 - Gold Section

3 The Gauss-Markov Theorem

- Statement
- Proof

4 Fisher Linear Discriminant

- History
- The Projection and The Rotation Idea
- Classifiers as Machines for dimensionality reduction
- Solution
 - Use the mean of each Class
 - Scatter measure
- The Cost Function
- A Transformation for simplification and defining the cost function
- **Where is this used?**
 - **Applications**
- Relation with Least Squared Error
 - What?

5 Exercises

- Some Stuff for you to try



Applications

Something Notable

- Bankruptcy prediction

- ▶ In bankruptcy prediction based on accounting ratios and other financial variables, linear discriminant analysis was the first statistical method applied to systematically explain which firms entered bankruptcy vs. survived.

- Face recognition

- ▶ In computerized face recognition, each face is represented by a large number of pixel values.
- ▶ The linear combinations obtained using Fisher's linear discriminant are called Fisher faces.

- Marketing

- ▶ In marketing, discriminant analysis was once often used to determine the factors which distinguish different types of customers and/or products on the basis of surveys or other forms of collected data.

Applications

Something Notable

- Bankruptcy prediction
 - ▶ In bankruptcy prediction based on accounting ratios and other financial variables, linear discriminant analysis was the first statistical method applied to systematically explain which firms entered bankruptcy vs. survived.
- Face recognition
 - ▶ In computerized face recognition, each face is represented by a large number of pixel values.
 - ▶ The linear combinations obtained using Fisher's linear discriminant are called Fisher faces.
- Marketing
 - ▶ In marketing, discriminant analysis was once often used to determine the factors which distinguish different types of customers and/or products on the basis of surveys or other forms of collected data.

Applications

Something Notable

- Bankruptcy prediction
 - ▶ In bankruptcy prediction based on accounting ratios and other financial variables, linear discriminant analysis was the first statistical method applied to systematically explain which firms entered bankruptcy vs. survived.
- Face recognition
 - ▶ In computerized face recognition, each face is represented by a large number of pixel values.
 - ▶ The linear combinations obtained using Fisher's linear discriminant are called Fisher faces.
- Marketing
 - ▶ In marketing, discriminant analysis was once often used to determine the factors which distinguish different types of customers and/or products on the basis of surveys or other forms of collected data.

Applications

Something Notable

- Bankruptcy prediction
 - ▶ In bankruptcy prediction based on accounting ratios and other financial variables, linear discriminant analysis was the first statistical method applied to systematically explain which firms entered bankruptcy vs. survived.
- Face recognition
 - ▶ In computerized face recognition, each face is represented by a large number of pixel values.
 - ▶ The linear combinations obtained using Fisher's linear discriminant are called Fisher faces.
- Marketing
 - ▶ In marketing, discriminant analysis was once often used to determine the factors which distinguish different types of customers and/or products on the basis of surveys or other forms of collected data.

Applications

Something Notable

- Bankruptcy prediction
 - ▶ In bankruptcy prediction based on accounting ratios and other financial variables, linear discriminant analysis was the first statistical method applied to systematically explain which firms entered bankruptcy vs. survived.
- Face recognition
 - ▶ In computerized face recognition, each face is represented by a large number of pixel values.
 - ▶ The linear combinations obtained using Fisher's linear discriminant are called Fisher faces.
- Marketing
 - ▶ In marketing, discriminant analysis was once often used to determine the factors which distinguish different types of customers and/or products on the basis of surveys or other forms of collected data.

Applications

Something Notable

- Bankruptcy prediction
 - ▶ In bankruptcy prediction based on accounting ratios and other financial variables, linear discriminant analysis was the first statistical method applied to systematically explain which firms entered bankruptcy vs. survived.
- Face recognition
 - ▶ In computerized face recognition, each face is represented by a large number of pixel values.
 - ▶ The linear combinations obtained using Fisher's linear discriminant are called Fisher faces.
- Marketing
 - ▶ In marketing, discriminant analysis was once often used to determine the factors which distinguish different types of customers and/or products on the basis of surveys or other forms of collected data.

Applications

Something Notable

- Bankruptcy prediction
 - ▶ In bankruptcy prediction based on accounting ratios and other financial variables, linear discriminant analysis was the first statistical method applied to systematically explain which firms entered bankruptcy vs. survived.
- Face recognition
 - ▶ In computerized face recognition, each face is represented by a large number of pixel values.
 - ▶ The linear combinations obtained using Fisher's linear discriminant are called Fisher faces.
- Marketing
 - ▶ In marketing, discriminant analysis was once often used to determine the factors which distinguish different types of customers and/or products on the basis of surveys or other forms of collected data.

Applications

Something Notable

- Bankruptcy prediction
 - ▶ In bankruptcy prediction based on accounting ratios and other financial variables, linear discriminant analysis was the first statistical method applied to systematically explain which firms entered bankruptcy vs. survived.
- Face recognition
 - ▶ In computerized face recognition, each face is represented by a large number of pixel values.
 - ▶ The linear combinations obtained using Fisher's linear discriminant are called Fisher faces.
- Marketing
 - ▶ In marketing, discriminant analysis was once often used to determine the factors which distinguish different types of customers and/or products on the basis of surveys or other forms of collected data.

Something Notable

- Biomedical studies

- ▶ The main application of discriminant analysis in medicine is the assessment of severity state of a patient and prognosis of disease outcome.



Something Notable

- Biomedical studies
 - ▶ The main application of discriminant analysis in medicine is the assessment of severity state of a patient and prognosis of disease outcome.



Please

Your Reading Material, it is about the Multi-class

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



Cinvestav

Outline

1 More in Regularization

- Introduction
- Smoothness of the Estimation
- The Error Estimate
- Choosing approximate inverses
- A Classic Example, Regularization as a Filter
- Another Example, The Landweber Iteration

2 Linear Regression using Gradient Descent

- Introduction
- What is the Gradient of the Equation?
- The Basic Algorithm
- How to obtain $\eta(k)$
 - Gold Section

3 The Gauss-Markov Theorem

- Statement
- Proof

4 Fisher Linear Discriminant

- History
- The Projection and The Rotation Idea
- Classifiers as Machines for dimensionality reduction
- Solution
 - Use the mean of each Class
 - Scatter measure
- The Cost Function
- A Transformation for simplification and defining the cost function
- Where is this used?
 - Applications
- Relation with Least Squared Error
 - What?

5 Exercises

- Some Stuff for you to try



Outline

1 More in Regularization

- Introduction
- Smoothness of the Estimation
- The Error Estimate
- Choosing approximate inverses
- A Classic Example, Regularization as a Filter
- Another Example, The Landweber Iteration

2 Linear Regression using Gradient Descent

- Introduction
- What is the Gradient of the Equation?
- The Basic Algorithm
- How to obtain $\eta(k)$
 - Gold Section

3 The Gauss-Markov Theorem

- Statement
- Proof

4 Fisher Linear Discriminant

- History
- The Projection and The Rotation Idea
- Classifiers as Machines for dimensionality reduction
- Solution
 - Use the mean of each Class
 - Scatter measure
- The Cost Function
- A Transformation for simplification and defining the cost function
- Where is this used?
 - Applications
- Relation with Least Squared Error
 - What?

5 Exercises

- Some Stuff for you to try



Relation with Least Squared Error

First

The least-squares approach to the determination of a linear discriminant was based on the goal of making the model predictions **as close as possible** to a set of target values.

Second

The Fisher criterion was derived by requiring maximum class separation in the output space.



Relation with Least Squared Error

First

The least-squares approach to the determination of a linear discriminant was based on the goal of making the model predictions **as close as possible** to a set of target values.

Second

The Fisher criterion was derived by requiring maximum class separation in the output space.



How do we do this?

First

- We have N samples.
- We have N_1 samples for class C_1 .
- We have N_2 samples for class C_2 .



How do we do this?

First

- We have N samples.
- We have N_1 samples for class C_1 .
- We have N_2 samples for class C_2 .

So, we decide the following for the targets on each class:

- We have then for class C_1 is $t_1 = \frac{N}{N_1}$.
- We have then for class C_2 is $t_2 = -\frac{N}{N_2}$.



How do we do this?

First

- We have N samples.
- We have N_1 samples for class C_1 .
- We have N_2 samples for class C_2 .

So, we decide the following for the targets for each class:

- We have then for class C_1 is $t_1 = \frac{N}{N_1}$.
- We have then for class C_2 is $t_2 = -\frac{N}{N_2}$.



How do we do this?

First

- We have N samples.
- We have N_1 samples for class C_1 .
- We have N_2 samples for class C_2 .

So, we decide the following for the targets on each class

- We have then for class C_1 is $t_1 = \frac{N}{N_1}$.
- We have then for class C_2 is $t_2 = -\frac{N}{N_2}$.



Thus

The new cost function (Our Classic Linear Model)

$$E = \frac{1}{2} \sum_{n=1}^N \left(\mathbf{w}^T \mathbf{x}_n + w_0 - t_n \right)^2 \quad (29)$$

Deriving with respect to w

$$\sum_{n=1}^N \left(\mathbf{w}^T \mathbf{x}_n + w_0 - t_n \right) \mathbf{x}_n = 0 \quad (30)$$

Deriving with respect to w_0

$$\sum_{n=1}^N \left(\mathbf{w}^T \mathbf{x}_n + w_0 - t_n \right) = 0 \quad (31)$$

Thus

The new cost function (Our Classic Linear Model)

$$E = \frac{1}{2} \sum_{n=1}^N (\mathbf{w}^T \mathbf{x}_n + w_0 - t_n)^2 \quad (29)$$

Deriving with respect to \mathbf{w}

$$\sum_{n=1}^N (\mathbf{w}^T \mathbf{x}_n + w_0 - t_n) \mathbf{x}_n = 0 \quad (30)$$

Deriving with respect to w_0

$$\sum_{n=1}^N (\mathbf{w}^T \mathbf{x}_n + w_0 - t_n) = 0 \quad (31)$$

Thus

The new cost function (Our Classic Linear Model)

$$E = \frac{1}{2} \sum_{n=1}^N (\mathbf{w}^T \mathbf{x}_n + w_0 - t_n)^2 \quad (29)$$

Deriving with respect to \mathbf{w}

$$\sum_{n=1}^N (\mathbf{w}^T \mathbf{x}_n + w_0 - t_n) \mathbf{x}_n = 0 \quad (30)$$

Deriving with respect to w_0

$$\sum_{n=1}^N (\mathbf{w}^T \mathbf{x}_n + w_0 - t_n) = 0 \quad (31)$$

Then

We have that

$$\sum_{n=1}^N (\mathbf{w}^T \mathbf{x}_n + w_0 - t_n) = \sum_{n=1}^N (\mathbf{w}^T \mathbf{x}_n + w_0) - \sum_{n=1}^N t_n$$

$$= \sum_{n=1}^N (\mathbf{w}^T \mathbf{x}_n + w_0) - N_1 \frac{N}{N_1} + N_2 \frac{N}{N_2}$$

$$= \sum_{n=1}^N (\mathbf{w}^T \mathbf{x}_n + w_0)$$

Then

We have that

$$\begin{aligned}\sum_{n=1}^N (\mathbf{w}^T \mathbf{x}_n + w_0 - t_n) &= \sum_{n=1}^N (\mathbf{w}^T \mathbf{x}_n + w_0) - \sum_{n=1}^N t_n \\ &= \sum_{n=1}^N (\mathbf{w}^T \mathbf{x}_n + w_0) - N_1 \frac{N}{N_1} + N_2 \frac{N}{N_2} \\ &= \sum_{n=1}^N (\mathbf{w}^T \mathbf{x}_n + w_0)\end{aligned}$$

Then

$$\left(\sum_{n=1}^N \mathbf{w}^T \mathbf{x}_n \right) + N w_0 = 0$$

Then

We have that

$$\begin{aligned}\sum_{n=1}^N (\mathbf{w}^T \mathbf{x}_n + w_0 - t_n) &= \sum_{n=1}^N (\mathbf{w}^T \mathbf{x}_n + w_0) - \sum_{n=1}^N t_n \\ &= \sum_{n=1}^N (\mathbf{w}^T \mathbf{x}_n + w_0) - N_1 \frac{N}{N_1} + N_2 \frac{N}{N_2} \\ &= \sum_{n=1}^N (\mathbf{w}^T \mathbf{x}_n + w_0)\end{aligned}$$

Then

$$\left(\sum_{n=1}^N \mathbf{w}^T \mathbf{x}_n \right) + N w_0 = 0$$

Then

We have that

$$w_0 = -\mathbf{w}^T \left(\frac{1}{N} \sum_{n=1}^N \mathbf{x}_n \right)$$

We rename $\frac{1}{N} \sum_{n=1}^N \mathbf{x}_n = \mathbf{m}$

$$\mathbf{m} = \frac{1}{N} \sum_{n=1}^N \mathbf{x}_n = \frac{1}{N} [N_1 \mathbf{m}_1 + N_2 \mathbf{m}_2]$$

Finally

$$w_0 = -\mathbf{w}^T \mathbf{m}$$



Then

We have that

$$w_0 = -\mathbf{w}^T \left(\frac{1}{N} \sum_{n=1}^N \mathbf{x}_n \right)$$

We rename $\frac{1}{N} \sum_{n=1}^N \mathbf{x}_n = \mathbf{m}$

$$\mathbf{m} = \frac{1}{N} \sum_{n=1}^N \mathbf{x}_n = \frac{1}{N} [N_1 \mathbf{m}_1 + N_2 \mathbf{m}_2]$$

Finally

$$w_0 = -\mathbf{w}^T \mathbf{m}$$



Then

We have that

$$w_0 = -\mathbf{w}^T \left(\frac{1}{N} \sum_{n=1}^N \mathbf{x}_n \right)$$

We rename $\frac{1}{N} \sum_{n=1}^N \mathbf{x}_n = \mathbf{m}$

$$\mathbf{m} = \frac{1}{N} \sum_{n=1}^N \mathbf{x}_n = \frac{1}{N} [N_1 \mathbf{m}_1 + N_2 \mathbf{m}_2]$$

Finally

$$w_0 = -\mathbf{w}^T \mathbf{m}$$



Now

In a similar way

$$\sum_{n=1}^N (\mathbf{w}^T \mathbf{x}_n + w_0) \mathbf{x}_n - \sum_{n=1}^N t_n \mathbf{x}_n = 0$$



Cinvestav

Thus, we have

Something Notable

$$\sum_{n=1}^N \left(\mathbf{w}^T \mathbf{x}_n + w_0 \right) \mathbf{x}_n - \frac{N}{N_1} \sum_{n=1}^{N_1} \mathbf{x}_n + \frac{N}{N_2} \sum_{n=1}^{N_2} \mathbf{x}_n = 0$$

Thus

$$\sum_{n=1}^N \left(\mathbf{w}^T \mathbf{x}_n + w_0 \right) \mathbf{x}_n - N (m_1 - m_2) = 0$$



Thus, we have

Something Notable

$$\sum_{n=1}^N (\mathbf{w}^T \mathbf{x}_n + w_0) \mathbf{x}_n - \frac{N}{N_1} \sum_{n=1}^{N_1} \mathbf{x}_n + \frac{N}{N_2} \sum_{n=1}^{N_2} \mathbf{x}_n = 0$$

Thus

$$\sum_{n=1}^N (\mathbf{w}^T \mathbf{x}_n + w_0) \mathbf{x}_n - N(m_1 - m_2) = 0$$



Next

Then, using $w_0 = -\mathbf{w}^T \mathbf{m}$

$$\sum_{n=1}^N (\mathbf{w}^T \mathbf{x}_n - \mathbf{w}^T \mathbf{m}) \mathbf{x}_n = N (\mathbf{m}_1 - \mathbf{m}_2)$$

Thus

$$\left[\sum_{n=1}^N (\mathbf{w}^T \mathbf{x}_n - \mathbf{w}^T \mathbf{m}) \mathbf{x}_n \right] = N (\mathbf{m}_1 - \mathbf{m}_2)$$



Next

Then, using $w_0 = -\mathbf{w}^T \mathbf{m}$

$$\sum_{n=1}^N (\mathbf{w}^T \mathbf{x}_n - \mathbf{w}^T \mathbf{m}) \mathbf{x}_n = N (\mathbf{m}_1 - \mathbf{m}_2)$$

Thus

$$\left[\sum_{n=1}^N (\mathbf{w}^T \mathbf{x}_n - \mathbf{w}^T \mathbf{m}) \mathbf{x}_n \right] = N (\mathbf{m}_1 - \mathbf{m}_2)$$



Now, Do you have the solution?

You have a version in Duda and Hart Section 5.8

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

This

$$\mathbf{X}^T \mathbf{X} \hat{\mathbf{w}} = \mathbf{X}^T \mathbf{y}$$

Now, we rewrite the data matrix

$$\mathbf{X} = \begin{bmatrix} \mathbf{1}_1 & \mathbf{X}_1 \\ -\mathbf{1}_2 & -\mathbf{X}_2 \end{bmatrix}$$



Now, Do you have the solution?

You have a version in Duda and Hart Section 5.8

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

Thus

$$\mathbf{X}^T \mathbf{X} \hat{\mathbf{w}} = \mathbf{X}^T \mathbf{y}$$

Now, we rewrite the data matrix

$$\mathbf{X} = \begin{bmatrix} 1_1 & X_1 \\ -1_2 & -X_2 \end{bmatrix}$$



Now, Do you have the solution?

You have a version in Duda and Hart Section 5.8

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

Thus

$$\mathbf{X}^T \mathbf{X} \hat{\mathbf{w}} = \mathbf{X}^T \mathbf{y}$$

Now, we rewrite the data matrix

$$\mathbf{X} = \begin{bmatrix} \mathbf{1}_1 & \mathbf{X}_1 \\ -\mathbf{1}_2 & -\mathbf{X}_2 \end{bmatrix}$$



In addition

Our old augmented w

$$w = \begin{bmatrix} w_0 \\ w \end{bmatrix}$$

And our new y

$$y = \begin{bmatrix} \frac{N}{N_1} 1_1 \\ \frac{N}{N_2} 1_2 \end{bmatrix} \quad (32)$$



In addition

Our old augmented w

$$w = \begin{bmatrix} w_0 \\ w \end{bmatrix}$$

And our new y

$$y = \begin{bmatrix} \frac{N}{N_1} \mathbf{1}_1 \\ \frac{N}{N_2} \mathbf{1}_2 \end{bmatrix} \quad (32)$$



Thus, we have

Something Notable

$$\begin{bmatrix} \mathbf{1}_1^T & -\mathbf{1}_2^T \\ \mathbf{X}_1^T & -\mathbf{X}_2^T \end{bmatrix} \begin{bmatrix} \mathbf{1}_1 & \mathbf{X}_1 \\ -\mathbf{1}_2 & -\mathbf{X}_2 \end{bmatrix} \begin{bmatrix} w_0 \\ \mathbf{w} \end{bmatrix} = \begin{bmatrix} \mathbf{1}_1^T & -\mathbf{1}_2^T \\ \mathbf{X}_1^T & -\mathbf{X}_2^T \end{bmatrix} \begin{bmatrix} \frac{N}{N_1} \mathbf{1}_1 \\ \frac{N}{N_2} \mathbf{1}_2 \end{bmatrix}$$



Thus, we have

Something Notable

$$\begin{bmatrix} \mathbf{1}_1^T & -\mathbf{1}_2^T \\ \mathbf{X}_1^T & -\mathbf{X}_2^T \end{bmatrix} \begin{bmatrix} \mathbf{1}_1 & \mathbf{X}_1 \\ -\mathbf{1}_2 & -\mathbf{X}_2 \end{bmatrix} \begin{bmatrix} w_0 \\ \mathbf{w} \end{bmatrix} = \begin{bmatrix} \mathbf{1}_1^T & -\mathbf{1}_2^T \\ \mathbf{X}_1^T & -\mathbf{X}_2^T \end{bmatrix} \begin{bmatrix} \frac{N}{N_1} \mathbf{1}_1 \\ \frac{N}{N_2} \mathbf{1}_2 \end{bmatrix}$$

Thus, if we use the following definitions for $i = 1, 2$

- $\mathbf{m}_i = \frac{1}{N_i} \sum_{\mathbf{x} \in C_i} \mathbf{x}$

- $S_w =$

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



Thus, we have

Something Notable

$$\begin{bmatrix} \mathbf{1}_1^T & -\mathbf{1}_2^T \\ \mathbf{X}_1^T & -\mathbf{X}_2^T \end{bmatrix} \begin{bmatrix} \mathbf{1}_1 & \mathbf{X}_1 \\ -\mathbf{1}_2 & -\mathbf{X}_2 \end{bmatrix} \begin{bmatrix} w_0 \\ \mathbf{w} \end{bmatrix} = \begin{bmatrix} \mathbf{1}_1^T & -\mathbf{1}_2^T \\ \mathbf{X}_1^T & -\mathbf{X}_2^T \end{bmatrix} \begin{bmatrix} \frac{N}{N_1} \mathbf{1}_1 \\ \frac{N}{N_2} \mathbf{1}_2 \end{bmatrix}$$

Thus, if we use the following definitions for $i = 1, 2$

- $\mathbf{m}_i = \frac{1}{N_i} \sum_{\mathbf{x} \in C_i} \mathbf{x}$
- $S_w = \sum_{\mathbf{x}_i \in C_1} (\mathbf{x}_i - \mathbf{m}_1) (\mathbf{x}_i - \mathbf{m}_1)^T + \sum_{\mathbf{x}_i \in C_2} (\mathbf{x}_i - \mathbf{m}_2) (\mathbf{x}_i - \mathbf{m}_2)^T$



Then

If we multiply the previous matrices

$$\begin{bmatrix} N & (N_1 \mathbf{m}_1 + N_2 \mathbf{m}_2)^T \\ (N_1 \mathbf{m}_1 + N_2 \mathbf{m}_2) & S_w + N_1 \mathbf{m}_1 \mathbf{m}_1^T + N_2 \mathbf{m}_2 \mathbf{m}_2^T \end{bmatrix} \begin{bmatrix} w_0 \\ \mathbf{w} \end{bmatrix} = \begin{bmatrix} 0 \\ N [\mathbf{m}_1 - \mathbf{m}_2] \end{bmatrix}$$

Then

$$\begin{bmatrix} N w_0 + (N_1 \mathbf{m}_1 + N_2 \mathbf{m}_2)^T \mathbf{w} \\ (N_1 \mathbf{m}_1 + N_2 \mathbf{m}_2) w_0 + [S_w + N_1 \mathbf{m}_1 \mathbf{m}_1^T + N_2 \mathbf{m}_2 \mathbf{m}_2^T] \mathbf{w} \end{bmatrix} = \begin{bmatrix} 0 \\ N [\mathbf{m}_1 - \mathbf{m}_2] \end{bmatrix}$$



Then

If we multiply the previous matrices

$$\begin{bmatrix} N & (N_1 \mathbf{m}_1 + N_2 \mathbf{m}_2)^T \\ (N_1 \mathbf{m}_1 + N_2 \mathbf{m}_2) & S_w + N_1 \mathbf{m}_1 \mathbf{m}_1^T + N_2 \mathbf{m}_2 \mathbf{m}_2^T \end{bmatrix} \begin{bmatrix} w_0 \\ \mathbf{w} \end{bmatrix} = \begin{bmatrix} 0 \\ N [\mathbf{m}_1 - \mathbf{m}_2] \end{bmatrix}$$

Then

$$\begin{bmatrix} N w_0 + (N_1 \mathbf{m}_1 + N_2 \mathbf{m}_2)^T \mathbf{w} \\ (N_1 \mathbf{m}_1 + N_2 \mathbf{m}_2) w_0 + [S_w + N_1 \mathbf{m}_1 \mathbf{m}_1^T + N_2 \mathbf{m}_2 \mathbf{m}_2^T] \mathbf{w} \end{bmatrix} = \begin{bmatrix} 0 \\ N [\mathbf{m}_1 - \mathbf{m}_2] \end{bmatrix}$$



Thus

We have that

$$\left[\frac{1}{N} S_w + \frac{N_1 N_2}{N^2} (\mathbf{m}_1 - \mathbf{m}_2) (\mathbf{m}_1 - \mathbf{m}_2)^T \right] \mathbf{w} = \mathbf{m}_1 - \mathbf{m}_2$$
$$w_0 = -\mathbf{w}^T \mathbf{m}$$

This

Since the vector $(\mathbf{m}_1 - \mathbf{m}_2) (\mathbf{m}_1 - \mathbf{m}_2)^T \mathbf{w}$ is in the direction of $\mathbf{m}_1 - \mathbf{m}_2$

$$\alpha = \frac{N_1 N_2}{N^2} (\mathbf{m}_1 - \mathbf{m}_2)^T \mathbf{w}$$

We have that

$$\frac{1}{N} S_w \mathbf{w} = (1 - \alpha) (\mathbf{m}_1 - \mathbf{m}_2)$$

Thus

We have that

$$\left[\frac{1}{N} S_w + \frac{N_1 N_2}{N^2} (\mathbf{m}_1 - \mathbf{m}_2) (\mathbf{m}_1 - \mathbf{m}_2)^T \right] \mathbf{w} = \mathbf{m}_1 - \mathbf{m}_2$$
$$w_0 = -\mathbf{w}^T \mathbf{m}$$

Thus

Since the vector $(\mathbf{m}_1 - \mathbf{m}_2) (\mathbf{m}_1 - \mathbf{m}_2)^T \mathbf{w}$ is in the direction of $\mathbf{m}_1 - \mathbf{m}_2$

$$\alpha = \frac{N_1 N_2}{N^2} (\mathbf{m}_1 - \mathbf{m}_2)^T \mathbf{w}$$

We have that

$$\frac{1}{N} S_w \mathbf{w} = (1 - \alpha) (\mathbf{m}_1 - \mathbf{m}_2)$$

Thus

We have that

$$w_0 = -\mathbf{w}^T \mathbf{m}$$
$$\left[\frac{1}{N} S_w + \frac{N_1 N_2}{N^2} (\mathbf{m}_1 - \mathbf{m}_2) (\mathbf{m}_1 - \mathbf{m}_2)^T \right] \mathbf{w} = \mathbf{m}_1 - \mathbf{m}_2$$

Thus

Since the vector $(\mathbf{m}_1 - \mathbf{m}_2) (\mathbf{m}_1 - \mathbf{m}_2)^T \mathbf{w}$ is in the direction of $\mathbf{m}_1 - \mathbf{m}_2$

$$\alpha = \frac{N_1 N_2}{N^2} (\mathbf{m}_1 - \mathbf{m}_2)^T \mathbf{w}$$

We have that

$$\frac{1}{N} S_w \mathbf{w} = (1 - \alpha) (\mathbf{m}_1 - \mathbf{m}_2)$$

Finally

We have that

$$\mathbf{w} = (1 - \alpha) N S_w^{-1} (\mathbf{m}_1 - \mathbf{m}_2) \propto S_w^{-1} (\mathbf{m}_1 - \mathbf{m}_2) \quad (33)$$



Cinvestav

Outline

1 More in Regularization

- Introduction
- Smoothness of the Estimation
- The Error Estimate
- Choosing approximate inverses
- A Classic Example, Regularization as a Filter
- Another Example, The Landweber Iteration

2 Linear Regression using Gradient Descent

- Introduction
- What is the Gradient of the Equation?
- The Basic Algorithm
- How to obtain $\eta(k)$
 - Gold Section

3 The Gauss-Markov Theorem

- Statement
- Proof

4 Fisher Linear Discriminant

- History
- The Projection and The Rotation Idea
- Classifiers as Machines for dimensionality reduction
- Solution
 - Use the mean of each Class
 - Scatter measure
- The Cost Function
- A Transformation for simplification and defining the cost function
- Where is this used?
 - Applications
- Relation with Least Squared Error
 - What?

5 Exercises

- Some Stuff for you to try



Machine Learning Theodoridis

Chapter 7

- 7.10, 7.13

Bishop

Chapter 4

- 4.4, 4.5, 4.6, 4.8



Machine Learning Theodoridis

Chapter 7

- 7.10, 7.13

Bishop

Chapter 4

- 4.4, 4.5, 4.6, 4.8

