

# Introduction to Machine Learning

## Gradient and Regularization Methods

Andres Mendez-Vazquez

February 10, 2023

# Outline

## 1 Linear Regression using Gradient Descent

- Introduction
- How do we stabilize the solution?
- The Basic Algorithm
- How to obtain  $\eta(k)$

## 2 Regularization Methods

- Introduction
- Intuition from Overfitting
- The Idea of Regularization
- Ridge Regression
- The LASSO
  - Lagrange Multipliers
  - The Basic Method
  - The Lagrangian Version of the LASSO
- Generalization

# Outline

## 1 Linear Regression using Gradient Descent

### ● Introduction

- How do we stabilize the solution?
- The Basic Algorithm
- How to obtain  $\eta(k)$

## 2 Regularization Methods

- Introduction
- Intuition from Overfitting
- The Idea of Regularization
- Ridge Regression
- The LASSO
  - Lagrange Multipliers
  - The Basic Method
  - The Lagrangian Version of the LASSO
- Generalization

# Given that the Canonical Solution has problems

We can develop a more robust algorithm

- Using the Gradient Descent Idea

# 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  $w(1)$ .

# Gradient Descent

The basic procedure is as follow

- 1 Start with a random weight vector  $\boldsymbol{w}(1)$ .
- 2 Compute the gradient vector  $\nabla J(\boldsymbol{w}(1))$ .

# 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:



# 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:

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

# Gradient Descent

The basic procedure is as follow

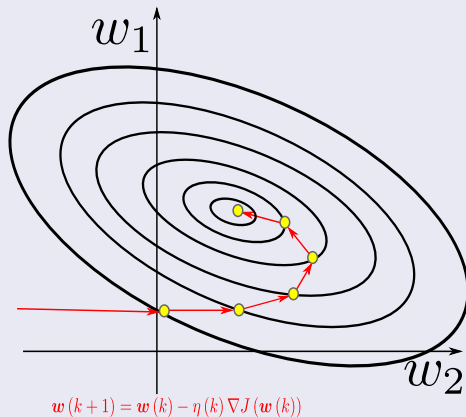
- 1 Start with a random weight vector  $\mathbf{w}(1)$ .
- 2 Compute the gradient vector  $\nabla J(\mathbf{w}(1))$ .
- 3 Obtain value  $\mathbf{w}(2)$  by moving from  $\mathbf{w}(1)$  in the direction of the steepest descent:

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

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

# Geometrically

We have the following



# Outline

## 1 Linear Regression using Gradient Descent

- Introduction
- **How do we stabilize the solution?**
- The Basic Algorithm
- How to obtain  $\eta(k)$

## 2 Regularization Methods

- Introduction
- Intuition from Overfitting
- The Idea of Regularization
- Ridge Regression
- The LASSO
  - Lagrange Multipliers
  - The Basic Method
  - The Lagrangian Version of the LASSO
- Generalization

## By using a regularized equation

We have

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

## By using a 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  $w_j$

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

## By using a 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 Linear Regression using Gradient Descent

- Introduction
- How do we stabilize the solution?
- **The Basic Algorithm**
- How to obtain  $\eta(k)$

## 2 Regularization Methods

- Introduction
- Intuition from Overfitting
- The Idea of Regularization
- Ridge Regression
- The LASSO
  - Lagrange Multipliers
  - The Basic Method
  - The Lagrangian Version of the LASSO
- Generalization



# Algorithm

## Gradient Decent

- 1 Initialize  $\mathbf{w}$ , criterion  $\theta$ ,  $\eta(\cdot)$ ,  $k = 0$

# Algorithm

## Gradient Decent

- 1 Initialize  $\mathbf{w}$ , criterion  $\theta$ ,  $\eta(\cdot)$ ,  $k = 0$
- 2 do  $k = k + 1$

# Algorithm

## Gradient Decent

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

# Algorithm

## Gradient Decent

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

# Algorithm

## Gradient Decent

- 1 Initialize  $\mathbf{w}$ , criterion  $\theta$ ,  $\eta(\cdot)$ ,  $k = 0$
- 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  $\mathbf{w}$

# Algorithm

## Gradient Decent

- 1 Initialize  $\mathbf{w}$ , criterion  $\theta$ ,  $\eta(\cdot)$ ,  $k = 0$
- 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!!!

# Algorithm

## Gradient Decent

- 1 Initialize  $\mathbf{w}$ , criterion  $\theta$ ,  $\eta(\cdot)$ ,  $k = 0$
- 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 Linear Regression using Gradient Descent

- Introduction
- How do we stabilize the solution?
- The Basic Algorithm
- How to obtain  $\eta(k)$

## 2 Regularization Methods

- Introduction
- Intuition from Overfitting
- The Idea of Regularization
- Ridge Regression
- The LASSO
  - Lagrange Multipliers
  - The Basic Method
  - The Lagrangian Version of the LASSO
- Generalization



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

We do the following

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

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 $\mathbf{w}(k)$

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

- $\nabla J$  is the vector of partial derivatives  $\frac{\partial J}{\partial w_i}$  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

- $\nabla J$  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)

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

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

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

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

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

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

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



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!!!

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 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 optimize 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 optimize using linear search methods

## Linear Search Methods

- Backtracking linear search

# 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 optimize using linear search methods

## Linear Search Methods

- Backtracking linear search
- Bisection method

# 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 optimize using linear search methods

## Linear Search Methods

- Backtracking linear search
- Bisection method
- Golden ratio

# 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 optimize using linear search methods

## Linear Search Methods

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



## Please Take a Look

For more, please read the paper

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

# Outline

## 1 Linear Regression using Gradient Descent

- Introduction
- How do we stabilize the solution?
- The Basic Algorithm
- How to obtain  $\eta(k)$

## 2 Regularization Methods

- **Introduction**
- Intuition from Overfitting
- The Idea of Regularization
- Ridge Regression
- The LASSO
  - Lagrange Multipliers
  - The Basic Method
  - The Lagrangian Version of the LASSO
- Generalization

# Shrinkage Methods

By retaining a subset of the predictors and discarding the rest

- Subset Selection produces a model that is interpretable,

# Shrinkage Methods

By retaining a subset of the predictors and discarding the rest

- Subset Selection produces a model that is interpretable,
- It possibly produces lower prediction error than the full model.

# Shrinkage Methods

By retaining a subset of the predictors and discarding the rest

- Subset Selection produces a model that is interpretable,
- It possibly produces lower prediction error than the full model.

However given process

- it often exhibits high variance,

# Shrinkage Methods

By retaining a subset of the predictors and discarding the rest

- Subset Selection produces a model that is interpretable,
- It possibly produces lower prediction error than the full model.

However given process

- it often exhibits high variance,
- It does not reduce the prediction error of the full model.

# Shrinkage Methods

By retaining a subset of the predictors and discarding the rest

- Subset Selection produces a model that is interpretable,
- It possibly produces lower prediction error than the full model.

However given process

- it often exhibits high variance,
- It does not reduce the prediction error of the full model.

Therefore

- Shrinkage methods are more continuous avoiding high variability.

# Outline

## 1 Linear Regression using Gradient Descent

- Introduction
- How do we stabilize the solution?
- The Basic Algorithm
- How to obtain  $\eta(k)$

## 2 Regularization Methods

- Introduction
- **Intuition from Overfitting**
- The Idea of Regularization
- Ridge Regression
- The LASSO
  - Lagrange Multipliers
  - The Basic Method
  - The Lagrangian Version of the LASSO
  - Generalization



# The house example

Imagine the following data set



Now assume that we use LSE

For the fitting

$$\frac{1}{2} \sum_{i=1}^N (h_{\mathbf{w}}(x_i) - y_i)^2$$

Now assume that we use LSE

For the fitting

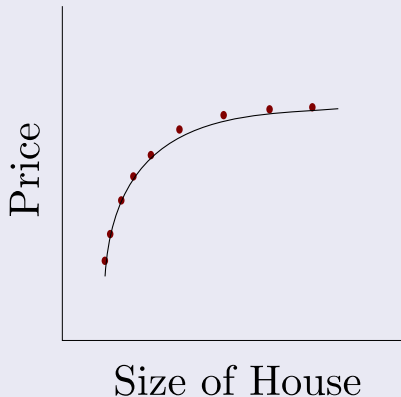
$$\frac{1}{2} \sum_{i=1}^N (h_{\mathbf{w}}(x_i) - y_i)^2$$

We can then run one of our machine to see what minimize better the previous equation

Question: Did you notice that I did not impose any structure to  $h_{\mathbf{w}}(x)$ ?

Then, First fitting

What about using  $h_1(x) = w_0 + w_1x + w_2x^2$ ?



## Second fitting

What about using  $h_2(x) = w_0 + w_1x + w_2x^2 + w_3x^3 + w_4x^4 + w_5x^5$ ?



Therefore, we have a problem

We get weird overfitting effects!!!

What do we do? What about minimizing the influence of  $w_3, w_4, w_5$ ?

Therefore, we have a problem

We get weird overfitting effects!!!

What do we do? What about minimizing the influence of  $w_3, w_4, w_5$ ?

How do we do that?

$$\min_w \frac{1}{2} \sum_{i=1}^N (h_w(x_i) - y_i)^2$$

What about integrating those values to the cost function? Ideas

# Outline

## 1 Linear Regression using Gradient Descent

- Introduction
- How do we stabilize the solution?
- The Basic Algorithm
- How to obtain  $\eta(k)$

## 2 Regularization Methods

- Introduction
- Intuition from Overfitting
- **The Idea of Regularization**
- Ridge Regression
- The LASSO
  - Lagrange Multipliers
  - The Basic Method
  - The Lagrangian Version of the LASSO
  - Generalization



We have

Regularization intuition is as follow

Small values for parameters  $w_0, w_1, w_2, \dots, w_n$

We have

Regularization intuition is as follow

Small values for parameters  $w_0, w_1, w_2, \dots, w_n$

It implies

- 1 "Simpler" function
- 2 Less prone to overfitting

We can do the previous idea for the other parameters

We can do the same for the other parameters

$$\min_w \frac{1}{2} \sum_{i=1}^N (h_w(x_i) - y_i)^2 + \sum_{i=1}^d \lambda_i w_i^2 \quad (9)$$

We can do the previous idea for the other parameters

We can do the same for the other parameters

$$\min_w \frac{1}{2} \sum_{i=1}^N (h_w(x_i) - y_i)^2 + \sum_{i=1}^d \lambda_i w_i^2 \quad (9)$$

However handling such many parameters can be so difficult

Combinatorial problem in reality!!!

Better, we can

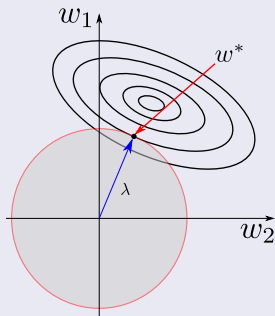
We better use the following

$$\min_{\mathbf{w}} \frac{1}{2} \sum_{i=1}^N (h_{\mathbf{w}}(x_i) - y_i)^2 + \lambda \sum_{i=1}^d w_i^2 \quad (10)$$

# Graphically

## Geometrically Equivalent to

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



# Outline

## 1 Linear Regression using Gradient Descent

- Introduction
- How do we stabilize the solution?
- The Basic Algorithm
- How to obtain  $\eta(k)$

## 2 Regularization Methods

- Introduction
- Intuition from Overfitting
- The Idea of Regularization
- **Ridge Regression**
- The LASSO
  - Lagrange Multipliers
  - The Basic Method
  - The Lagrangian Version of the LASSO
  - Generalization

# Ridge Regression

## Equation

$$\hat{\mathbf{w}} = \arg \min_{\mathbf{w}} \left\{ \sum_{i=1}^N \left( y_i - w_0 - \sum_{j=1}^d x_{ij} w_j \right)^2 + \lambda \sum_{j=1}^d w_j^2 \right\}$$



# Ridge Regression

## Equation

$$\hat{\mathbf{w}} = \arg \min_{\mathbf{w}} \left\{ \sum_{i=1}^N \left( y_i - w_0 - \sum_{j=1}^d x_{ij} w_j \right)^2 + \lambda \sum_{j=1}^d w_j^2 \right\}$$

## Here

- $\lambda \geq 0$  is a complexity parameter that controls the amount of shrinkage

Therefore

The Larger  $\lambda \geq 0$

- The coefficients are shrunk toward zero (and each other).

# Therefore

The Larger  $\lambda \geq 0$

- The coefficients are shrunk toward zero (and each other).

This is also used in Neural Networks

- where it is known as weight decay

This is also can be written

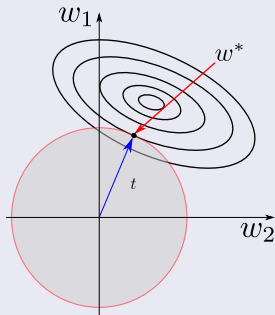
### Optimization Solution

$$\begin{aligned} & \arg \min_{\mathbf{w}} \sum_{i=1}^N \left( y_i - w_0 - \sum_{j=1}^d x_{ij} w_j \right)^2 \\ & \text{subject to } \sum_{j=1}^d w_j^2 < t \end{aligned}$$

# Graphically

## Geometrically Equivalent to

$$\begin{aligned} \arg \min \quad & \sum_{i=1}^N (y_i - \mathbf{x}_i^T \mathbf{w})^2 \\ \text{subject to} \quad & \sum_{i=1}^{d+1} w_i^2 < t \end{aligned}$$



# Outline

## 1 Linear Regression using Gradient Descent

- Introduction
- How do we stabilize the solution?
- The Basic Algorithm
- How to obtain  $\eta(k)$

## 2 Regularization Methods

- Introduction
- Intuition from Overfitting
- The Idea of Regularization
- Ridge Regression
- **The LASSO**
  - Lagrange Multipliers
  - The Basic Method
  - The Lagrangian Version of the LASSO
  - Generalization

# Least Absolute Shrinkage and Selection Operator (LASSO)

It was introduced by Robert Tibshirani in 1996 based on Leo Breiman's nonnegative garrote

$$\hat{\mathbf{w}}^{garrote} = \arg \min_{\mathbf{w}} \sum_{i=1}^N \left( y_i - \beta_0 - \sum_{j=1}^d x_{ij} w_j \right)^2 + N\lambda \sum_{j=1}^d w_j$$

s.t.  $w_j > 0 \ \forall j$

# Least Absolute Shrinkage and Selection Operator (LASSO)

It was introduced by Robert Tibshirani in 1996 based on Leo Breiman's nonnegative garrote

$$\hat{\mathbf{w}}^{garrote} = \arg \min_{\mathbf{w}} \sum_{i=1}^N \left( y_i - \beta_0 - \sum_{j=1}^d x_{ij} w_j \right)^2 + N\lambda \sum_{j=1}^d w_j$$

s.t.  $w_j > 0 \ \forall j$

This is quite derivable

However, Tibshirani realized that you could get a more flexible model by using the absolute value at the constraint!!!



# Least Absolute Shrinkage and Selection Operator (LASSO)

It was introduced by Robert Tibshirani in 1996 based on Leo Breiman's nonnegative garrote

$$\hat{\mathbf{w}}^{garrote} = \arg \min_{\mathbf{w}} \sum_{i=1}^N \left( y_i - \beta_0 - \sum_{j=1}^d x_{ij} w_j \right)^2 + N\lambda \sum_{j=1}^d w_j$$

s.t.  $w_j > 0 \ \forall j$

This is quite derivable

However, Tibshirani realized that you could get a more flexible model by using the absolute value at the constraint!!!

Robert Tibshirani proposed the use of the  $L_1$  norm

$$\|\mathbf{w}\|_1 = \sum_{i=1}^d |w_i|$$

# The Final Optimization Problem

## LASSO

$$\begin{aligned}\hat{\mathbf{w}}^{LASSO} &= \arg \min_{\mathbf{w}} \sum_{i=1}^N \left( y_i - \beta_0 - \sum_{j=1}^d x_{ij} w_j \right)^2 \\ \text{s.t. } &\sum_{i=1}^d |w_i| \leq t\end{aligned}$$

# The Final Optimization Problem

## LASSO

$$\begin{aligned}\hat{\mathbf{w}}^{LASSO} = \arg \min_{\mathbf{w}} \sum_{i=1}^N \left( y_i - \beta_0 - \sum_{j=1}^d x_{ij} w_j \right)^2 \\ \text{s.t. } \sum_{i=1}^d |w_i| \leq t\end{aligned}$$

This is not derivable

- More advanced methods are necessary to solve this problem!!!

# Outline

## 1 Linear Regression using Gradient Descent

- Introduction
- How do we stabilize the solution?
- The Basic Algorithm
- How to obtain  $\eta(k)$

## 2 Regularization Methods

- Introduction
- Intuition from Overfitting
- The Idea of Regularization
- Ridge Regression
- **The LASSO**
  - **Lagrange Multipliers**
  - The Basic Method
  - The Lagrangian Version of the LASSO
  - Generalization

# Lagrange Multipliers

## The method of Lagrange multipliers

- It gives a set of necessary conditions to identify optimal points of equality constrained optimization problems.

# Lagrange Multipliers

## The method of Lagrange multipliers

- It gives a set of necessary conditions to identify optimal points of equality constrained optimization problems.

This is done by converting a constrained problem to an equivalent unconstrained problem

- with the help of certain unspecified parameters known as Lagrange multipliers.

# Lagrange Multipliers

## The classical problem formulation

$$\begin{aligned} \min \quad & f(x_1, \dots, x_n) \\ \text{s.t.} \quad & h_1(x_1, \dots, x_n) = 0 \end{aligned}$$

# Lagrange Multipliers

## The classical problem formulation

$$\begin{aligned} \min \quad & f(x_1, \dots, x_n) \\ \text{s.t.} \quad & h_1(x_1, \dots, x_n) = 0 \end{aligned}$$

## It can be converted into

$$\min L(x_1, \dots, x_n, \lambda) = \min \{f(x_1, \dots, x_n) - \lambda h_1(x_1, \dots, x_n)\}$$



# Lagrange Multipliers

## The classical problem formulation

$$\begin{aligned} \min \quad & f(x_1, \dots, x_n) \\ \text{s.t.} \quad & h_1(x_1, \dots, x_n) = 0 \end{aligned}$$

## It can be converted into

$$\min L(x_1, \dots, x_n, \lambda) = \min \{f(x_1, \dots, x_n) - \lambda h_1(x_1, \dots, x_n)\}$$

## where

- $L(\mathbf{x}, \lambda)$  is the Lagrangian function.
- $\lambda$  is an unspecified positive or negative constant called the **Lagrange Multiplier**.

# Finding an Optimum using Lagrange Multipliers

## New problem

$$\min L(x_1, \dots, x_n, \lambda) = \min \{f(x_1, \dots, x_n) - \lambda h_1(x_1, \dots, x_n)\}$$

# Finding an Optimum using Lagrange Multipliers

## New problem

$$\min L(x_1, \dots, x_n, \lambda) = \min \{f(x_1, \dots, x_n) - \lambda h_1(x_1, \dots, x_n)\}$$

We want a  $\lambda = \lambda^*$  optimal

If the minimum of  $L(x_1, \dots, x_n, \lambda^*)$  occurs at

$$(x_1, x_2, \dots, x_n)^T = (x_1, x_2, \dots, x_n)^{T*}$$

Therefore

$(x_1, \dots, x_n)^{T*}$  satisfies  $h_1(x_1, \dots, x_n) = 0$ , then  $(x_1, \dots, x_n)^{T*}$  minimizes

$$\begin{aligned} \min \quad & f(x_1, \dots, x_n) \\ \text{s.t.} \quad & h_1(x_1, \dots, x_n) = 0 \end{aligned}$$

Therefore

$(x_1, \dots, x_n)^{T*}$  satisfies  $h_1(x_1, \dots, x_n) = 0$ , then  $(x_1, \dots, x_n)^{T*}$  minimizes

$$\begin{aligned} \min f(x_1, \dots, x_n) \\ \text{s.t. } h_1(x_1, \dots, x_n) = 0 \end{aligned}$$

Trick

- It is to find appropriate value for Lagrangian multiplier  $\lambda$ .

# Remember

Think about this

Remember First Law of Newton!!!

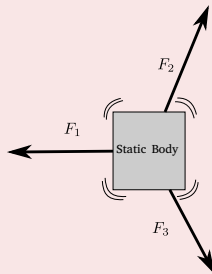
# Remember

Think about this

Remember First Law of Newton!!!

Yes!!!

**A system in equilibrium does not move**



# Lagrange Multipliers

## Definition

Gives a set of necessary conditions to identify optimal points of equality constrained optimization problem



# Lagrange was a Physicists

He was thinking in the following formula

A system in equilibrium has the following equation:

$$F_1 + F_2 + \dots + F_K = 0 \quad (11)$$

# Lagrange was a Physicists

He was thinking in the following formula

A system in equilibrium has the following equation:

$$F_1 + F_2 + \dots + F_K = 0 \quad (11)$$

But functions do not have forces?

Are you sure?

# Lagrange was a Physicists

He was thinking in the following formula

A system in equilibrium has the following equation:

$$F_1 + F_2 + \dots + F_K = 0 \quad (11)$$

But functions do not have forces?

Are you sure?

Think about the following

The Gradient of a surface.

# Gradient to a Surface

After all a gradient is a measure of the maximal change

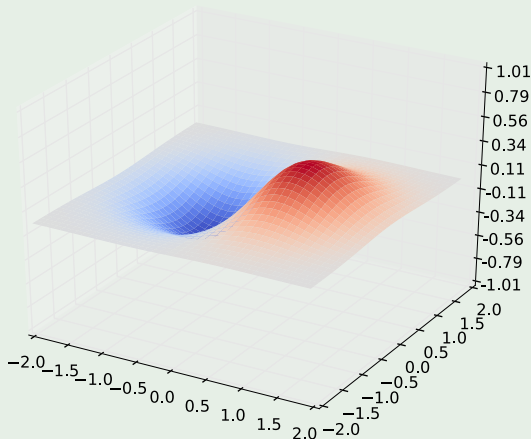
For example the gradient of a function of three variables:

$$\nabla f(\mathbf{x}) = i \frac{\partial f(\mathbf{x})}{\partial x} + j \frac{\partial f(\mathbf{x})}{\partial y} + k \frac{\partial f(\mathbf{x})}{\partial z} \quad (12)$$

where  $i$ ,  $j$  and  $k$  are unitary vectors in the directions  $x$ ,  $y$  and  $z$ .

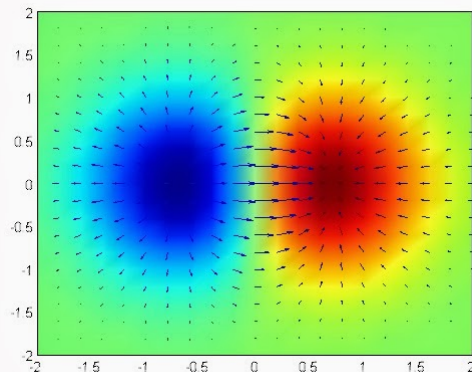
## Example

We have  $f(x, y) = x \exp \{-x^2 - y^2\}$



# Example

With Gradient at the the contours when projecting in the 2D plane



100

Now, Think about this

Yes, we can use the gradient

However, we need to do some scaling of the forces by using parameters  $\lambda$

## Now, Think about this

Yes, we can use the gradient

However, we need to do some scaling of the forces by using parameters  $\lambda$

Thus, we have

$$F_0 + \lambda_1 F_1 + \dots + \lambda_K F_K = 0 \quad (13)$$

where  $F_0$  is the gradient of the principal cost function and  $F_i$  for  $i = 1, 2, \dots, K$ .



## Thus

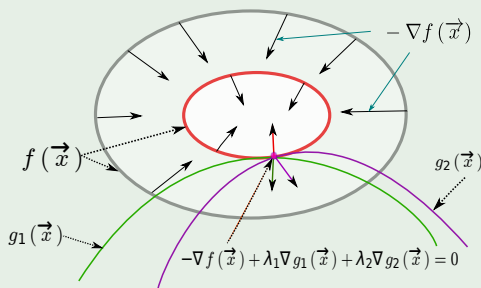
If we have the following optimization:

$$\begin{aligned} \min \quad & f(\mathbf{x}) \\ \text{s.t.} \quad & g_1(\mathbf{x}) = 0 \\ & g_2(\mathbf{x}) = 0 \end{aligned}$$

# Geometric interpretation in the case of minimization

What is wrong? Gradients are going in the other direction, we can fix by simple multiplying by -1

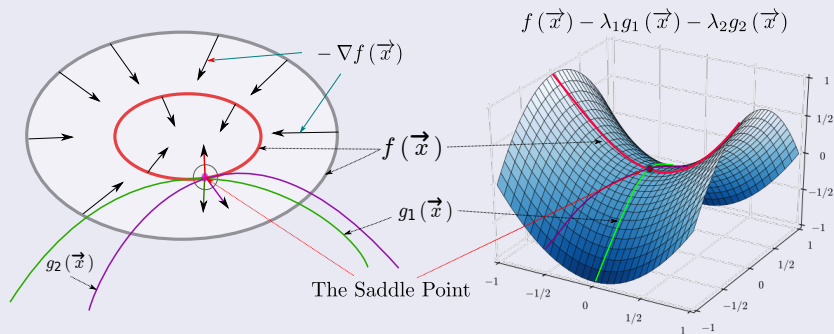
Here the cost function is  $f(x, y) = x \exp \{-x^2 - y^2\}$  we want to minimize



Nevertheless: it is equivalent to  $\nabla f(\vec{x}) - \lambda_1 \nabla g_1(\vec{x}) - \lambda_2 \nabla g_2(\vec{x}) = 0$

Basically, we convert the problem into a one looking for a **Saddle Point**

At the left the original problem, at the right the Lagrangian!!!



# Yes!!!

## Basically

- We convert the minimization or maximization of a convex or concave section of a function living in a constrained environment!!!

# Outline

## 1 Linear Regression using Gradient Descent

- Introduction
- How do we stabilize the solution?
- The Basic Algorithm
- How to obtain  $\eta(k)$

## 2 Regularization Methods

- Introduction
- Intuition from Overfitting
- The Idea of Regularization
- Ridge Regression
- **The LASSO**
  - Lagrange Multipliers
  - **The Basic Method**
  - The Lagrangian Version of the LASSO
  - Generalization

# Method

## Steps

- 1 Original problem is rewritten as:
  - 1 minimize  $L(\mathbf{x}, \lambda) = f(\mathbf{x}) - \lambda h_1(\mathbf{x})$

# Method

## Steps

- 1 Original problem is rewritten as:
  - 1 minimize  $L(\mathbf{x}, \lambda) = f(\mathbf{x}) - \lambda h_1(\mathbf{x})$
- 2 Take derivatives of  $L(\mathbf{x}, \lambda)$  with respect to  $x_i$  and set them equal to zero.

# Method

## Steps

- 1 Original problem is rewritten as:
  - 1 minimize  $L(\mathbf{x}, \lambda) = f(\mathbf{x}) - \lambda h_1(\mathbf{x})$
- 2 Take derivatives of  $L(\mathbf{x}, \lambda)$  with respect to  $x_i$  and set them equal to zero.

$$\sum_{i=1}^N \left( y_i - \mathbf{x}^T \mathbf{w} \right)^2 + \lambda \sum_{i=1}^d |w_i| \quad (14)$$

- 1 Here you need to use a soft version of the absolute value
- 3 Express all  $x_i$  in terms of Lagrangian multiplier  $\lambda$ .



# Method

## Steps

- 1 Original problem is rewritten as:
  - 1 minimize  $L(\mathbf{x}, \lambda) = f(\mathbf{x}) - \lambda h_1(\mathbf{x})$
- 2 Take derivatives of  $L(\mathbf{x}, \lambda)$  with respect to  $x_i$  and set them equal to zero.

$$\sum_{i=1}^N \left( y_i - \mathbf{x}^T \mathbf{w} \right)^2 + \lambda \sum_{i=1}^d |w_i| \quad (14)$$

- 1 Here you need to use a soft version of the absolute value
- 3 Express all  $x_i$  in terms of Lagrangian multiplier  $\lambda$ .
- 4 Plug  $\mathbf{x}$  in terms of  $\lambda$  in constraint  $h_1(\mathbf{x}) = 0$  and solve  $\lambda$ .

# Method

## Steps

- 1 Original problem is rewritten as:
  - 1 minimize  $L(\mathbf{x}, \lambda) = f(\mathbf{x}) - \lambda h_1(\mathbf{x})$
- 2 Take derivatives of  $L(\mathbf{x}, \lambda)$  with respect to  $x_i$  and set them equal to zero.

$$\sum_{i=1}^N \left( y_i - \mathbf{x}^T \mathbf{w} \right)^2 + \lambda \sum_{i=1}^d |w_i| \quad (14)$$

- 1 Here you need to use a soft version of the absolute value
- 3 Express all  $x_i$  in terms of Lagrangian multiplier  $\lambda$ .
- 4 Plug  $\mathbf{x}$  in terms of  $\lambda$  in constraint  $h_1(\mathbf{x}) = 0$  and solve  $\lambda$ .
- 5 Calculate  $\mathbf{x}$  by using the just found value for  $\lambda$ .

# Method

## Steps

- 1 Original problem is rewritten as:
  - 1 minimize  $L(\mathbf{x}, \lambda) = f(\mathbf{x}) - \lambda h_1(\mathbf{x})$
- 2 Take derivatives of  $L(\mathbf{x}, \lambda)$  with respect to  $x_i$  and set them equal to zero.

$$\sum_{i=1}^N \left( y_i - \mathbf{x}^T \mathbf{w} \right)^2 + \lambda \sum_{i=1}^d |w_i| \quad (14)$$

- 1 Here you need to use a soft version of the absolute value
- 3 Express all  $x_i$  in terms of Lagrangian multiplier  $\lambda$ .
- 4 Plug  $\mathbf{x}$  in terms of  $\lambda$  in constraint  $h_1(\mathbf{x}) = 0$  and solve  $\lambda$ .
- 5 Calculate  $\mathbf{x}$  by using the just found value for  $\lambda$ .

## From the step 2

If there are  $n$  variables (i.e.,  $x_1, \dots, x_n$ ) then you will get  $n$  equations with  $n + 1$  unknowns (i.e.,  $n$  variables  $x_i$  and one Lagrangian multiplier  $\lambda$ ).

## Example

We can apply that to the following problem

$$\begin{aligned} \min \quad & f(x, y) = x^2 - 8x + y^2 - 12y + 48 \\ \text{s.t.} \quad & x + y = 8 \end{aligned}$$

# Outline

## 1 Linear Regression using Gradient Descent

- Introduction
- How do we stabilize the solution?
- The Basic Algorithm
- How to obtain  $\eta(k)$

## 2 Regularization Methods

- Introduction
- Intuition from Overfitting
- The Idea of Regularization
- Ridge Regression
- **The LASSO**
  - Lagrange Multipliers
  - The Basic Method
  - **The Lagrangian Version of the LASSO**
  - Generalization

# The Lagrangian Version

## The Lagrangian

$$\hat{\mathbf{w}}^{LASSO} = \arg \min_{\mathbf{w}} \left\{ \sum_{i=1}^N \left( y_i - \mathbf{x}^T \mathbf{w} \right)^2 + \lambda \sum_{i=1}^d |w_i| \right\}$$

# The Lagrangian Version

## The Lagrangian

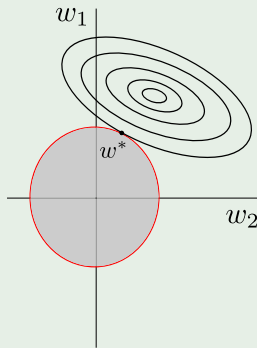
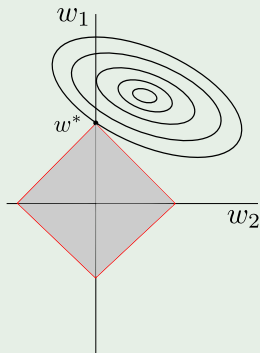
$$\hat{\mathbf{w}}^{LASSO} = \arg \min_{\mathbf{w}} \left\{ \sum_{i=1}^N \left( y_i - \mathbf{x}^T \mathbf{w} \right)^2 + \lambda \sum_{i=1}^d |w_i| \right\}$$

## However

- You have other regularizations as  $\|\mathbf{w}\|_2 = \sqrt{\sum_{i=1}^d |w_i|^2}$

# Graphically

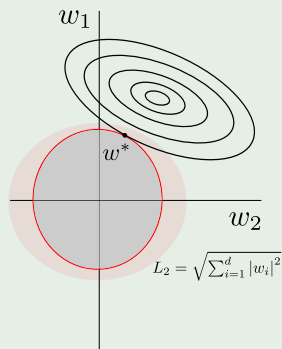
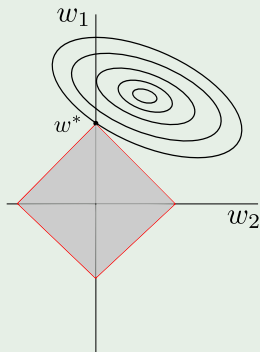
The first area correspond to the  $L_1$  regularization and the second one?





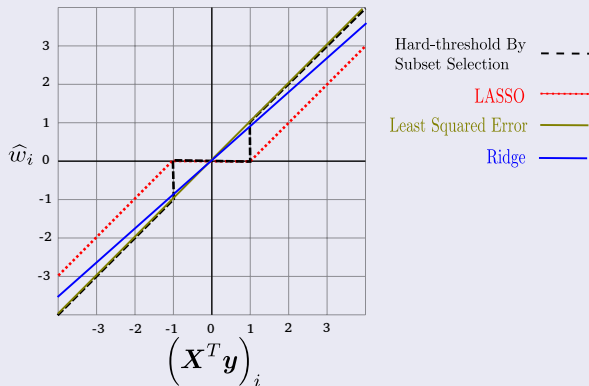
# Graphically

Yes the circle defined as  $\|w\|_2 = \sqrt{\sum_{i=1}^d |w_i|^2}$



# For Example

In the Case of  $X$  is a Orthogonal Matrix



# The seminal paper by Robert Tibshirani

An initial study of this regularization can be seen in

“Regression Shrinkage and Selection via the LASSO” by Robert Tibshirani  
- 1996

This out the scope of this class

However, it is worth noticing that the most efficient method for solving LASSO problems is

“Pathwise Coordinate Optimization” By Jerome Friedman, Trevor Hastie, Holger Ho and Robert Tibshirani

This out the scope of this class

However, it is worth noticing that the most efficient method for solving LASSO problems is

“Pathwise Coordinate Optimization” By Jerome Friedman, Trevor Hastie, Holger Ho and Robert Tibshirani

Nevertheless

It will be a great seminar paper!!!

# Outline

## 1 Linear Regression using Gradient Descent

- Introduction
- How do we stabilize the solution?
- The Basic Algorithm
- How to obtain  $\eta(k)$

## 2 Regularization Methods

- Introduction
- Intuition from Overfitting
- The Idea of Regularization
- Ridge Regression
- **The LASSO**
  - Lagrange Multipliers
  - The Basic Method
  - The Lagrangian Version of the LASSO
- **Generalization**

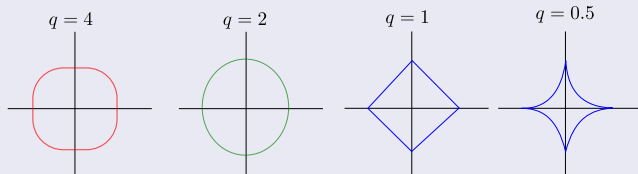
## Furthermore

We can generalize ridge regression and the lasso, and view them as Bayes estimates

$$\hat{\mathbf{w}}^{LASSO} = \arg \min_{\mathbf{w}} \left\{ \sum_{i=1}^N \left( y_i - \mathbf{x}^T \mathbf{w} \right)^2 + \lambda \sum_{i=1}^d |w_i|^q \right\} \text{ with } q \geq 0$$

## For Example

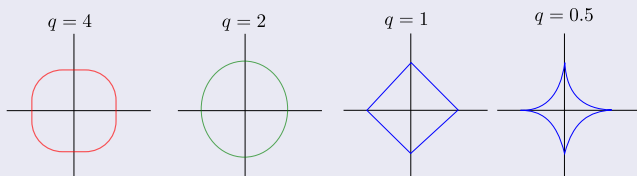
We have when  $d = 2$





## For Example

We have when  $d = 2$



Here, when  $q > 1$

- You are having a derivable Lagrangian, but you lose the LASSO properties

Therefore

Zou and Hastie (2005) introduced the elastic- net penalty

$$\lambda \sum_{i=1}^d \left\{ \alpha w_i^2 + (1 - \alpha) |w_i| \right\}$$

Therefore

Zou and Hastie (2005) introduced the elastic- net penalty

$$\lambda \sum_{i=1}^d \left\{ \alpha w_i^2 + (1 - \alpha) |w_i| \right\}$$

This is Basically

- A Compromise Between the Ridge and LASSO.