



# Optimization and Solving Non-linear Equations

MODULE DES130: COMPUTATIONAL STATISTICS

---

Dr. Erick A. Chacón Montalván (echacon@uni.edu.pe)

Escuela de Profesional de Ingeniería Estadística  
Facultad de Ingeniería Económica, Estadística y Ciencias Sociales  
Universidad Nacional de Ingeniería (UNI)  
Lima – Perú

# Introduction

---

# Optimization

Let the function  $f : A \rightarrow \mathcal{R}$ .

## Minimization:

Find  $x_0 \in A / f(x_0) \leq f(x)$  for all  $x \in A$ .

## Maximization:

Find  $x_0 \in A / f(x_0) \geq f(x)$  for all  $x \in A$ .

## Notes:

- Maximizing  $f(x)$  is equivalent to minimizing  $-f(x)$ .
- Local and global optimizers.
- It is not easy to distinguish local and global optimizers.
- Subfields of optimization depend on the properties of  $f(x)$

# Solving non-linear equations

Optimization is intimately linked with solving non-linear equations (Givens and Hoeting 2012).

## **Optimization as root-finding problem:**

Let the function  $f : A \rightarrow \mathcal{R}$ . Then, finding

$$x_0 \in A / f(x_0) \leq f(x) \text{ for all } x \in A.$$

Is equivalent to finding

$$x_0 \in A / f'(x_0) = 0 \text{ and } f''(x_0) > 0.$$

## Common optimization types

### **Linear programming:**

The function  $f(x)$  and restrictions are linear. Example, estimation of the coefficients of a quantile regression.

### **Quadratic programming:**

The function  $f(x)$  can be a quadratic function. Example, least squares estimation.

### **Non-linear programming:**

The function  $f(x)$  is non-linear. Example, maximum likelihood estimation.

# Quadratic programming

---

## Motivation: least squares

Let consider the following linear model  $y_i = \mathbf{x}_i^\top \boldsymbol{\beta} + e_i$  for  $i = 1, \dots, n$ , then the least squares estimation is done by the optimization

$$\hat{\boldsymbol{\beta}} = \arg \min_{\boldsymbol{\beta}} \sum_{i=1}^n (y_i - \mathbf{x}_i^\top \boldsymbol{\beta})^2.$$

It can be shown that  $\hat{\boldsymbol{\beta}} = (\mathbf{X}^\top \mathbf{X})^{-1} \mathbf{X}^\top \mathbf{y}$ . However, this expression is not really computed in practice (Wood 2017).

**Consider the QR decomposition:**

$$\mathbf{X}_{n \times p} = \mathbf{Q}_{n \times n} \begin{bmatrix} \mathbf{R}_{p \times p} \\ \mathbf{0} \end{bmatrix} = \mathbf{Q}_f \mathbf{R}, \quad \text{and} \quad \mathbf{Q}^\top \mathbf{y} = \begin{bmatrix} \mathbf{f}_{p \times 1} \\ \mathbf{r} \end{bmatrix},$$

where  $\mathbf{R}$  is a upper triangular matrix and  $\mathbf{Q}$  is an orthogonal matrix.

Using the matricial form of the loss function:

$$\begin{aligned}L(\beta) &= \|\mathbf{y} - \mathbf{X}\beta\|^2 \\&= \|\mathbf{Q}^T\mathbf{y} - \mathbf{Q}^T\mathbf{X}\beta\|^2 \\&= \left\| \mathbf{Q}^T\mathbf{y} - \begin{bmatrix} \mathbf{R} \\ \mathbf{0} \end{bmatrix} \beta \right\|^2 \\&= \left\| \begin{bmatrix} \mathbf{f} \\ \mathbf{r} \end{bmatrix} - \begin{bmatrix} \mathbf{R}\beta \\ \mathbf{0} \end{bmatrix} \right\|^2 \\L(\beta) &= \|\mathbf{f} - \mathbf{R}\beta\|^2 + \|\mathbf{r}\|^2.\end{aligned}$$

Which is minimized by  $\hat{\beta} = \mathbf{R}^{-1}\mathbf{f}$ . A similar algorithm can be found for generalized least squares (Wood 2017).



# Linear programming

---

Introduction

Quadratic programming

**Linear programming**

**Motivation: quantile regression**

Quantile regression as a linear program

Non-linear programming

References

## Motivation: quantile regression

Let  $Y_i$  with  $Pr(Y_i \leq Q_p) = p$  such as the conditional quantile is

$$Q_p(Y_i | \mathbf{x}_i) = \mathbf{x}_i^T \beta_p.$$

Estimation for the  $p$ -th quantile regression, defined by , is done by the optimization

$$\hat{\beta}_p = \arg \min_{\beta_p} \left( \sum_{i: y_i \geq \mathbf{x}_i^T \beta_p} p |y_i - \mathbf{x}_i^T \beta_p| + \sum_{i: y_i < \mathbf{x}_i^T \beta_p} (1 - p) |y_i - \mathbf{x}_i^T \beta_p| \right).$$

Introduction

Quadratic programming

**Linear programming**

Motivation: quantile regression

**Quantile regression as a linear program**

Non-linear programming

References

## Quantile regression as a linear program

Given  $\mathbf{y} - \mathbf{X}\beta_p = \mathbf{e} = \mathbf{r}^+ - \mathbf{r}^-$ , where  $r_i^+ = \max(0, e_i)$  and  $r_i^- = \max(0, -e_i)$ . Then the estimation of quantile regression can be expressed as a linear program:

### Linear programming problem:

$$\min_{\mathbf{r}^+, \mathbf{r}^-, \beta_p} (p\mathbf{1}^T \mathbf{r}^+ + (1-p)\mathbf{1}^T \mathbf{r}^-)$$

subject to

$$\mathbf{y} = \mathbf{X}\beta_p + \mathbf{r}^+ - \mathbf{r}^-,$$

$$\mathbf{r}^+ \in \mathbb{R}_+^n, \text{ and } \mathbf{r}^- \in \mathbb{R}_+^n.$$

# Non-linear programming

---

Introduction

Quadratic programming

Linear programming

**Non-linear programming**

**Introduction**

Newton's Method

Steepest ascent

Gauss-Newton Method

References

## Motivation: generalized linear models

### Exponential family:

$$f_{Y_i}(y_i) = \exp \left\{ \frac{y_i \theta - b(\theta)}{a(\phi)} + c(y, \phi) \right\}$$

Let consider the following model:

$$Y_i \mid \mathbf{x}_i \sim \text{Bernoulli}(\pi_i),$$

$$\text{logit}(\pi_i) = \mathbf{x}_i^\top \beta.$$

Given that  $f(y_i) = \exp \left\{ y_i \log \left( \frac{\pi_i}{1-\pi_i} \right) + \log(1 - \pi) \right\}$ , then

$$\theta_i = \text{logit}(\pi_i) = \mathbf{x}_i^\top \beta,$$

$$b(\theta) = \log(1 - \pi_i) = \log(1 + \exp(\theta_i)),$$

$$a(\phi) = 1, \text{ and } c(y_i, \phi) = 0.$$



## Taylor series approximation

It can be seen as a approximation of a function  $f : \mathbb{R} \rightarrow \mathbb{R}$  by a polinomial function  $g(x)$  around the point  $x_0$ .

**For**  $r = 0 \rightarrow g(x) = c_0$ , \ **given**  $g(x_0) = f(x_0)$ , \ **then**  $g(x) = f(x_0)$ .

**For**  $r = 1 \rightarrow g(x) = c_0 + c_1x$ , \ **given**  $g(x_0) = f(x_0)$ , and  $g'(x_0) = f'(x_0)$ , \ **then**  $g(x) = f(x_0) + f'(x_0)(x - x_0)$ .

**For**  $r = 2 \rightarrow g(x) = c_0 + c_1x + c_2x^2$ , \ **given**  $g(x_0) = f(x_0)$ ,  $g'(x_0) = f'(x_0)$ , and  $g''(x_0) = f''(x_0)$ , \ **then**  $g(x) = f(x_0) + f'(x_0)(x - x_0) + \frac{f''(x_0)}{2}(x - x_0)^2$ .

# Taylor series approximation

## In general:

Let  $f(x)$  be  $n$ -th differentiable, then the  $n$ -th order Taylor approximation of  $f(x)$  around  $x_0$  is

$$g(x) = \sum_{r=0}^n \frac{f^{(r)}(x_0)}{r!} (x - x_0)^r.$$

## Example:

Obtain the first order Taylor approximation of  $f(x) = \exp(x)$  around 0.

### Solución:

Notice that  $x_0 = 0$  such that  $f(0) = \exp(0) = 1$   $f'(0) = \exp(0) = 1$ .  
Hence  $g(x) = 1 + 1(x - 0) = 1 + x$ .

# Quadratic Taylor series approximation

## In general:

Let  $f : \mathbb{R}^n \rightarrow \mathbb{R}$  be a function taking as input  $\mathbf{x} \in \mathbb{R}^n$  with second partial derivatives, then the quadratic Taylor approximation of  $f(\mathbf{x})$  around  $\mathbf{x}_0$  is

$$g(\mathbf{x}) = f(\mathbf{x}_0) + (\mathbf{x} - \mathbf{x}_0)^\top \mathbf{f}'(\mathbf{x}_0) + \frac{1}{2}(\mathbf{x} - \mathbf{x}_0)^\top \mathbf{f}''(\mathbf{x}_0)(\mathbf{x} - \mathbf{x}_0).$$

Introduction

Quadratic programming

Linear programming

**Non-linear programming**

Introduction

**Newton's Method**

Steepest ascent

Gauss-Newton Method

References

## Newton's Method: definition

Also referred as *Newton-Raphson iteration*. The solution to  $f(x) = 0$  can be obtained by performing an approximation, and finding the solution to that approximation (Gentle 2009).

**Step 1:** Approximate  $f(x)$  by the first-order Taylor expansion.

$$f(x) \approx f(x^{(t)}) + f'(x^{(t)})(x - x^{(t)})$$

**Step 2:** Find the root of the approximation.

$$x^{(t+1)} = x^{(t)} - \frac{f(x^{(t)})}{f'(x^{(t)})}$$

## Newton's Method for optimization

The solution to  $f'(x) = 0$  can be obtained by performing an approximation, and finding the solution to that approximation.

**Step 1:** Approximate  $f'(x)$  by the first-order Taylor expansion.

$$f'(x) \approx f'(x^{(t)}) + f''(x^{(t)})(x - x^{(t)})$$

**Step 2:** Find the root of the approximation.

$$x^{(t+1)} = x^{(t)} - \frac{f'(x^{(t)})}{f''(x^{(t)})}$$

Using the Fisher information  $I(\theta) = \mathbb{E} [-f''(\theta)]$  instead of  $-f''(x^{(t)})$  leads to the Fisher scoring method.

## Newton's Method for optimization

More generally, for  $f : \mathbb{R}^n \rightarrow \mathbb{R}$ :

**Step 1:** Approximate  $f'(\mathbf{x})$  by the first-order Taylor expansion.

$$f'(\mathbf{x}) \approx f'(\mathbf{x}^{(t)}) + f''(\mathbf{x}^{(t)})(\mathbf{x} - \mathbf{x}^{(t)})$$

**Step 2:** Find the root of the approximation.

$$\mathbf{x}^{(t+1)} = \mathbf{x}^{(t)} - f''(\mathbf{x}^{(t)})^{-1}f'(\mathbf{x}^{(t)}).$$

Using the Fisher information  $\mathbf{I}(\theta) = \mathbb{E} [-f''(\theta)]$  instead of  $-f''(\mathbf{x}^{(t)})$  leads to the Fisher scoring method.

# GLM: Iteratively reweighted least squares

Let consider the following model:

$$Y_i | \mathbf{x}_i \sim \text{Bernoulli}(\pi_i),$$
$$\text{logit}(\pi_i) = \mathbf{x}_i^\top \boldsymbol{\beta}.$$

Where

$$\theta_i = \text{logit}(\pi_i) = \mathbf{x}_i^\top \boldsymbol{\beta},$$
$$b(\theta) = \log(1 - \pi_i) = \log(1 + \exp(\theta_i)) = \log(1 + \exp(\mathbf{x}_i^\top \boldsymbol{\beta})),$$
$$a(\phi) = 1, \text{ and } c(y_i, \phi) = 0.$$

Hence

## Likelihood function

$$l(\boldsymbol{\beta}) = \mathbf{y}^\top \mathbf{X} \boldsymbol{\beta} - \mathbf{b}^\top \mathbf{1}$$



## GLM: Iteratively reweighted least squares

Obtaining the gradient and Hessian:

$$l'(\beta) = \mathbf{X}^T(\mathbf{y} - \pi)$$

$$l''(\beta) = \frac{d}{d\beta}(\mathbf{X}^T(\mathbf{y} - \pi)) = -\mathbf{X}^T\mathbf{W}\mathbf{X}.$$

Where  $\mathbf{W}$  is diagonal with  $w_{ii} = \pi_i(1 - \pi_i)$ . Hence, applying Newton's method:

$$\begin{aligned}\beta^{(t+1)} &= \beta^{(t)} - l''(\beta^{(t)})^{-1}l'(\beta^{(t)}) \\ &= \beta^{(t)} - \left(\mathbf{X}^T\mathbf{W}^{(t)}\mathbf{X}\right)^{-1} \left(\mathbf{X}^T(\mathbf{y} - \pi^{(t)})\right). \\ &= \left(\mathbf{X}^T\mathbf{W}^{(t)}\mathbf{X}\right)^{-1}\mathbf{X}^T\mathbf{W}^{(t)}\mathbf{z}^{(t)}.\end{aligned}$$

where  $\mathbf{z}^{(t)} = \mathbf{X}\beta^{(t)} + (\mathbf{W}^{(t)})^{-1}(\mathbf{y} - \pi^{(t)})$ .

Introduction

Quadratic programming

Linear programming

**Non-linear programming**

Introduction

Newton's Method

**Steepest ascent**

Gauss-Newton Method

References

## Newton-like method

Update such as

$$\mathbf{x}^{(t+1)} = \mathbf{x}^{(t)} - (\mathbf{M}^{(t)})^{-1}g'(\mathbf{x}^t),$$

where  $\mathbf{M}^{(t)}$  is a Hessian's approximation  
[@givens2012computational].

## Steepest ascent

Replaces  $\mathbf{M}^{(t)}$  by  $-\mathbf{I}$ .

$$\mathbf{x}^{(t+1)} = \mathbf{x}^{(t)} + \alpha^{(t)}g'(\mathbf{x}^t),$$

Introduction

Quadratic programming

Linear programming

**Non-linear programming**

Introduction

Newton's Method

Steepest ascent

**Gauss-Newton Method**

References

Consider the following model

$$Y_i = h(\mathbf{x}_i, \boldsymbol{\theta}) + \varepsilon_i.$$

where  $h(\mathbf{x}_i, \boldsymbol{\theta})$  is non-linear. This model has loss function:

$$L(\boldsymbol{\theta}) = \|\mathbf{y} - \mathbf{h}(\mathbf{X}, \boldsymbol{\theta})\|^2.$$

Instead of approximating  $L(\boldsymbol{\theta})$ , we can approximate the function  $h(\mathbf{x}_i, \boldsymbol{\theta})$  using Taylor series.

## Gauss-Newton method

By approximating  $h(\mathbf{x}_i, \theta)$  we obtain

$$Y_i = h(\mathbf{x}_i, \theta^{(t)}) + (\theta - \theta^{(t)})^\top h'(\mathbf{x}_i, \theta^{(t)}) + \varepsilon_i.$$

$$Y_i = \tilde{h}(\mathbf{x}_i, \theta^{(t)}, \theta) + \varepsilon_i.$$

Let  $z_i^{(t)} = y_i - h(\mathbf{x}_i, \theta^{(t)})$ , and  $\mathbf{a}_i^{(t)} = h'(\mathbf{x}_i, \theta^{(t)})$ . Then

$$\mathbf{z}^{(t)} = \mathbf{A}^{(t)}(\theta - \theta^{(t)}) + \varepsilon.$$

The estimate can be obtained as follows:

$$\theta^{(t+1)} - \theta^{(t)} = \left( \mathbf{A}^{(t)\top} \mathbf{A}^{(t)} \right)^{-1} \mathbf{A}^{(t)\top} \mathbf{z}^{(t)}.$$

$$\theta^{(t+1)} = \theta^{(t)} + \left( \mathbf{A}^{(t)\top} \mathbf{A}^{(t)} \right)^{-1} \mathbf{A}^{(t)\top} \mathbf{z}^{(t)}.$$

It does not require to compute the Hessian matrix.

# References

---

- Gentle, James E. 2009. *Computational Statistics*. Springer Science & Business Media.
- Givens, Geof H., and Jennifer A. Hoeting. 2012. *Computational Statistics*. John Wiley & Sons.
- Wood, Simon N. 2017. *Generalized Additive Models: An Introduction with R, Second Edition*. CRC Press.