

# STA410 | Programming Portfolio Assignment 3

## 3.0 - Nonlinearity s| $\mathbf{b} = \mathbf{f}(\mathbf{x}) \neq \mathbf{A}\mathbf{x}$

A more interesting problem than  $\mathbf{Ax} = \mathbf{b}$  is the problem  $\mathbf{g}(\mathbf{x}) = \mathbf{b}$  for a **nonlinear function**  $\mathbf{g}$ . Such **nonlinear equations** are a general version form of a wide set of problems. However, a most frequently encountered instance of the problem is

$$f'(z) = 0 \quad \text{or more generally} \quad \underbrace{\nabla_z f(z)} = \mathbf{0} \quad \text{(multivariate form)}$$

because solutions to these equations are **local minima** or **maxima** of  $f$  and **model fitting** can be framed as a subclass of this general **optimization problem**.

Two important notes about these **optimization problems** can be made.

1. The **derivative**  $f'(z)$  is just some other function, say  $g(z)$ , so outside of its useful interpretation as a derivative,  $f'(z)$  may be treated just as any other function might be treated.

E.g., the **first-order Taylor series approximation** of  $f$  is simply

$$f(x) \approx f(x_0) + (x - x_0)f'(x_0)$$

as if  $f$  had been replaced with some other function  $g$ .

2. **Optimization solutions**  $f'(x^*) = 0$  are found within regions of curvature of  $f$ , while **roots**  $f(x_0) = 0$  need not be.
- The behavior of a function near a **solution** to an **optimization** problem differs, generally speaking, from its behavior near a **root**.
  - And the **numerical precision** in an **optimization** context to in general differs from the vanilla **root-finding** context.

E.g., to the degree that a **second order Taylor series approximation**  $g_{\tilde{x}}(x)$  of  $f(x)$  is accurate

$$f(x) \approx g_{\tilde{x}}(x) = f(\tilde{x}) + (x - \tilde{x})f'(\tilde{x}) + \frac{f''(\tilde{x})}{2}(x - \tilde{x})^2$$

changes in  $g_{\tilde{x}}(x)$  are likely dominated by

- the linear term  $\underbrace{(x - \tilde{x})}_{\epsilon} f'(\tilde{x}) \rightarrow 0$  if  $\tilde{x}$  is near **root**  $x_0$
- the quadratic term  $\frac{1}{2} \underbrace{(x - \tilde{x})^2}_{\epsilon^2} f''(\tilde{x}) \rightarrow 0$  if  $\tilde{x}$  is near **optimum**  $x^*$ , since  $f'(\tilde{x}) \approx 0$ .

I.e.,  $f(x) \approx g_{x^*}(x)$  evaluated near an **optimization problem solution**  $x^*$  only supports about half of the **numerical accuracy**

i.e., the square root of the available precision  $\epsilon_{\text{machine}}$ , since the difference will be squared

compared to evaluating  $f(x) \approx g_{x_0}(x)$  near one of its **roots**  $x_0$ .

\newpage

### 3.0.1 Maximum Likelihood Estimates (MLEs)

The **score function** is the gradient of the **log likelihood**

$$\nabla_{\theta} l(\theta) = \left( \frac{\partial l(\theta)}{\partial \theta_1}, \dots, \frac{\partial l(\theta)}{\partial \theta_p} \right)^T \quad \text{where} \quad l(\theta) = \log f(x|\theta) \overset{\text{iid}}{=} \log \prod_{i=1}^n f(x_i|\theta)$$

**Maximum Likelihood Estimates** (MLEs) come from the (**nonlinear**) **score equation** which sets the **score function** equal to  $\mathbf{0}$ . And for the **true value** of the parameter  $\theta^{\text{true}}$ , the **score function** has expected value  $\mathbf{0}$  (with respect to  $f_x$  the distribution of the data). So

$$\underbrace{\nabla_{\theta} l(\hat{\theta}) = \mathbf{0}}_{\text{score equation}} \quad \text{and} \quad \underbrace{E_X[\nabla_{\theta} l(\theta^{\text{true}})] = \mathbf{0}}_{\text{score function}}$$

The expected value of the **score function** follows since

$$\begin{aligned} E[\nabla_{\theta} l(\theta)] &= \int \nabla_{\theta} l(\theta) f(x|\theta) dx = \int \left( \frac{\partial l(\theta)}{\partial \theta_1}, \dots, \frac{\partial l(\theta)}{\partial \theta_p} \right)^T f(x|\theta) dx \\ &= \int \left( \frac{\partial}{\partial \theta_1} \log f(x|\theta), \dots, \frac{\partial}{\partial \theta_p} \log f(x|\theta) \right)^T f(x|\theta) dx \\ &= \int \left( \frac{1}{f(x|\theta)} \frac{\partial f(x|\theta)}{\partial \theta_1}, \dots, \frac{1}{f(x|\theta)} \frac{\partial f(x|\theta)}{\partial \theta_p} \right)^T f(x|\theta) dx \\ &= \int \left( \frac{\partial f(x|\theta)}{\partial \theta_1}, \dots, \frac{\partial f(x|\theta)}{\partial \theta_p} \right)^T dx = \nabla_{\theta} \int f(x|\theta) dx = \nabla_{\theta} 1 = \mathbf{0} \end{aligned}$$

The **Fisher information matrix**  $I(\theta^{\text{true}})$ , or **expected Fisher information matrix** is the expected value of the **outer product** of the **score function** with itself and is equal to the expected value of the negative of the **Hessian** of the log likelihood (all with respect to  $f_x$  the distribution of the data), i.e.,  $I(\theta^{\text{true}}) = E_X[\nabla_{\theta} l(\theta^{\text{true}}) \nabla_{\theta} l(\theta^{\text{true}})^T] = E_X[-H_{ll}(\theta^{\text{true}})]$

The **observed Fisher information** is

$$\begin{aligned} \hat{I}(\hat{\theta}) &= -H_{ll}(\hat{\theta}) = \underbrace{-J(\nabla_{\theta} l)(\hat{\theta})}_{\text{Jacobian of score function evaluated at } \hat{\theta}} \\ &\approx -\sum_{i=1}^n J(\nabla_{\theta} \log f(x_i|\theta)) \big|_{\hat{\theta}} \\ \hat{I}(\hat{\theta}) &\approx \sum_{i=1}^n \nabla_{\theta} \log f(x_i|\theta) \big|_{\hat{\theta}} \left( \nabla_{\theta} \log f(x_i|\theta) \big|_{\hat{\theta}} \right)^T \end{aligned}$$

And the **asymptotic distribution** of the MLE is

$$p(\hat{\theta}) \overset{n \rightarrow \infty}{\longrightarrow} N(\theta^{\text{true}}, \Sigma) \quad \Sigma = \frac{I(\theta^{\text{true}})^{-1}}{n} \approx \frac{\hat{I}(\hat{\theta})^{-1}}{n}$$

$$\mathbf{I}(\hat{\theta})^{-1} \approx \frac{1}{n} \mathbf{I}(\theta^{\text{true}})^{-1}$$

where either the **observed Fisher information** or its approximation based on the **outer product** of the **score function** with itself may be used as plug in estimates for the **expected Fisher information matrix**  $\mathbf{I}(\theta^{\text{true}})$ .