

1 Foundation

1.1 Matrix Derivatives

- Let $\mathbf{x} = (x_1, \dots, x_n)^T$ and f a real-valued (multivariate) function of \mathbf{x} . The derivative of $f(\mathbf{x})$, if exists, is given by

$$\frac{\partial f(\mathbf{x})}{\partial \mathbf{x}} = \begin{pmatrix} \frac{\partial f(\mathbf{x})}{\partial x_1} \\ \vdots \\ \frac{\partial f(\mathbf{x})}{\partial x_n} \end{pmatrix}$$

$$\frac{\partial f(\mathbf{x})}{\partial \mathbf{x}^T} = \left(\frac{\partial f(\mathbf{x})}{\partial \mathbf{x}} \right)^T$$

- The second-order partial derivative of f at \mathbf{x} is given by

$$\frac{\partial^2 f(\mathbf{x})}{\partial \mathbf{x} \partial \mathbf{x}^T}^{n \times n} = \begin{pmatrix} \frac{\partial^2 f(\mathbf{x})}{\partial^2 x_1} & \frac{\partial^2 f(\mathbf{x})}{\partial x_1 \partial x_2} & \cdots & \frac{\partial^2 f(\mathbf{x})}{\partial x_1 \partial x_n} \\ \frac{\partial^2 f(\mathbf{x})}{\partial x_2 \partial x_1} & \frac{\partial^2 f(\mathbf{x})}{\partial^2 x_2} & \cdots & \frac{\partial^2 f(\mathbf{x})}{\partial x_2 \partial x_n} \\ \vdots & \vdots & \cdots & \vdots \\ \frac{\partial^2 f(\mathbf{x})}{\partial x_n \partial x_1} & \frac{\partial^2 f(\mathbf{x})}{\partial x_n \partial x_2} & \cdots & \frac{\partial^2 f(\mathbf{x})}{\partial^2 x_n} \end{pmatrix}$$

- Let \mathbf{f} be a $p \times 1$ vector of functions of \mathbf{x} as follows,

$$\mathbf{f}(\mathbf{x}) = \begin{pmatrix} f_1(\mathbf{x}) \\ \vdots \\ f_p(\mathbf{x}) \end{pmatrix}$$

Then the derivative of \mathbf{f} at \mathbf{x} , if exists, is given by

$$\frac{\partial \mathbf{f}(\mathbf{x})}{\partial \mathbf{x}}^{n \times p} = \begin{pmatrix} \frac{\partial f_1(\mathbf{x})}{\partial x_1} & \frac{\partial f_2(\mathbf{x})}{\partial x_1} & \cdots & \frac{\partial f_p(\mathbf{x})}{\partial x_1} \\ \frac{\partial f_1(\mathbf{x})}{\partial x_2} & \frac{\partial f_2(\mathbf{x})}{\partial x_2} & \cdots & \frac{\partial f_p(\mathbf{x})}{\partial x_2} \\ \vdots & \vdots & \cdots & \vdots \\ \frac{\partial f_1(\mathbf{x})}{\partial x_n} & \frac{\partial f_2(\mathbf{x})}{\partial x_n} & \cdots & \frac{\partial f_p(\mathbf{x})}{\partial x_n} \end{pmatrix}$$

$$\frac{\partial \mathbf{f}(\mathbf{x})}{\partial \mathbf{x}^T} = \left(\frac{\partial \mathbf{f}(\mathbf{x})}{\partial \mathbf{x}} \right)^T$$

- Let $f(\mathbf{x}) = \mathbf{a}^T \mathbf{x} = \mathbf{x}^T \mathbf{a}$, then

$$\frac{\partial f(\mathbf{x})}{\partial \mathbf{x}} = \mathbf{a}$$

- Let $g(\mathbf{x}) = \mathbf{x}^T \mathbf{A} \mathbf{x}$, then

$$\frac{\partial g(\mathbf{x})}{\partial \mathbf{x}} = \mathbf{A} \mathbf{x} + (\mathbf{x}^T \mathbf{A})^T = \mathbf{A} \mathbf{x} + \mathbf{A}^T \mathbf{x}$$

- Let $\mathbf{f}(\mathbf{x}) = \mathbf{A}^T \mathbf{x}$, then

$$\frac{\partial \mathbf{f}(\mathbf{x})}{\partial \mathbf{x}} = \mathbf{A}$$

- Let $g(\mathbf{x}) = \mathbf{x}^T \mathbf{B}$, then

$$\frac{\partial g(\mathbf{x})}{\partial \mathbf{x}} = \mathbf{B}$$

- Let $\mathbf{A}(\mathbf{x}) = (A_{ij}(\mathbf{x}))_{n \times n}$ which is invertible, then

$$\frac{\partial \mathbf{A}^{-1}}{\partial x_i} = -\mathbf{A}^{-1} \frac{\partial \mathbf{A}}{\partial x_i} \mathbf{A}^{-1}$$

$$\frac{\partial |\mathbf{A}|}{\partial x_i} = \text{tr} \left(|\mathbf{A}| \mathbf{A}^{-1} \frac{\partial \mathbf{A}}{\partial x_i} \right) = |\mathbf{A}| \text{tr} \left(\mathbf{A}^{-1} \frac{\partial \mathbf{A}}{\partial x_i} \right)$$

1.2 Multivariate normal distribution

$\mathbf{Y} = (Y_1, Y_2, \dots, Y_n)^T \sim MVN_n(\boldsymbol{\mu}, \boldsymbol{\Sigma})$ for a positive definite $\boldsymbol{\Sigma}$ if any of the following definitions hold

- the density function of \mathbf{Y} is

$$f_{\mathbf{Y}}(\mathbf{y}) = \frac{1}{(2\pi)^{n/2} |\boldsymbol{\Sigma}|^{1/2}} \exp \left\{ -\frac{1}{2} (\mathbf{y} - \boldsymbol{\mu})^T \boldsymbol{\Sigma}^{-1} (\mathbf{y} - \boldsymbol{\mu}) \right\}$$

- the moment generating function of \mathbf{Y} is

$$m_{\mathbf{Y}}(\mathbf{t}) \equiv E \left(e^{\mathbf{t}^T \mathbf{Y}} \right) = \exp \left\{ \boldsymbol{\mu}^T \mathbf{t} + \frac{1}{2} \mathbf{t}^T \boldsymbol{\Sigma} \mathbf{t} \right\}$$

- \mathbf{Y} has the same distribution as

$$\mathbf{Q}^T \mathbf{z} + \boldsymbol{\mu},$$

where $\mathbf{z} = (z_1, z_2, \dots, z_n)^T$, z_i 's are iid $N(0, 1)$, $\mathbf{Q}^T \mathbf{Q} = \boldsymbol{\Sigma}$.

Property 1. $E(\mathbf{Y}) = \boldsymbol{\mu}$ and $\text{Var}(\mathbf{Y}) = \boldsymbol{\Sigma}$

Property 2. If z_i 's are iid standard normal $N(0, 1)$, then

$$\mathbf{z} = (z_1, z_2, \dots, z_n) \sim \text{MVN}(\mathbf{0}_n, \mathbf{I}_{n \times n})$$

Property 3. If $\mathbf{Y} \sim MVN_n(\boldsymbol{\mu}, \boldsymbol{\Sigma})$ and $\mathbf{A}_{p \times n}$ is not random, then

$$\mathbf{A}\mathbf{Y} \sim MVN_p(\mathbf{A}\boldsymbol{\mu}, \mathbf{A}\boldsymbol{\Sigma}\mathbf{A}^T)$$

Property 4. If $\mathbf{Y} \sim MVN_n(\boldsymbol{\mu}, \sigma^2 \mathbf{I}_{n \times n})$ and $\mathbf{A}_{n \times n}$ is a constant orthogonal matrix, then

$$\mathbf{A}\mathbf{Y} \sim MVN_n(\mathbf{A}\boldsymbol{\mu}, \sigma^2 \mathbf{I}_{n \times n})$$

1.3 Other relevant distributions

Chi-square Distribution with p degrees of freedom, $\chi^2(p)$: suppose Z_1, Z_2, \dots, Z_p are iid $N(0, 1)$, then

$$\sum_{i=1}^p Z_i^2 \sim \chi^2(p)$$

t-Distribution with p degrees of freedom, $t(p)$: suppose $Z \sim N(0, 1)$ and $V \sim \chi^2(p)$ and Z and V are independent, then

$$\frac{Z}{\sqrt{V/p}} \sim t(p)$$

F-Distribution with p and q degrees of freedom, $F(p, q)$: suppose $Z \sim \chi^2(p)$ and $V \sim \chi^2(q)$ and Z and V are independent, then

$$\frac{Z/p}{V/q} \sim F(p, q).$$

1.4 Exponential family

Exponential family is of the following form:

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

where θ is the *canonical parameter*, ϕ is the *dispersion parameter*, $b(\theta)$ is the *cumulative function* (different from cumulative generating function) and $a()$, $b()$, $c()$ are known functions.

Exponential family also has **another** commonly used form [A good material here](#)

$$p(x | \eta) = h(x) \exp \{ \eta^T T(x) - A(\eta) \}$$

parameter vector η , often referred to as the *canonical parameter* for given functions T and h . The statistic $T(X)$ is referred to as a *sufficient statistic*. The function $A(\eta)$ is known as the *cumulant function* and

$$A(\eta) = \log \int h(x) \exp \{ \eta^T T(x) \} \nu(dx)$$

with respect to the measure ν

Some examples

- $Y \sim N(\mu, \sigma^2)$

$$f(y; \mu, \sigma^2) = \frac{1}{\sqrt{2\pi\sigma^2}} \exp \left\{ -\frac{(y - \mu)^2}{2\sigma^2} \right\} = \exp \left\{ \frac{y\mu - \mu^2/2}{\sigma^2} - \frac{y^2}{2\sigma^2} - \frac{1}{2} \log(2\pi\sigma^2) \right\}$$

- $\theta = \mu, \phi = \sigma^2$
- $a(\phi) = \phi$
- $b(\theta) = \frac{\theta^2}{2}$
- $c(y, \phi) = -\frac{y^2}{2\phi} - \frac{1}{2} \log(2\pi\phi)$

with another formulation

$$p(x | \mu, \sigma^2) = \frac{1}{\sqrt{2\pi\sigma}} \exp \left\{ -\frac{1}{2\sigma^2}(x - \mu)^2 \right\} = \frac{1}{\sqrt{2\pi}} \exp \left\{ \frac{\mu}{\sigma^2}x - \frac{1}{2\sigma^2}x^2 - \frac{1}{2\sigma^2}\mu^2 - \log \sigma \right\}$$

then

$$\begin{aligned} \eta &= \begin{bmatrix} \mu/\sigma^2 \\ -1/2\sigma^2 \end{bmatrix} \\ T(x) &= \begin{bmatrix} x \\ x^2 \end{bmatrix} \\ A(\eta) &= \frac{\mu^2}{2\sigma^2} + \log \sigma = -\frac{\eta_1^2}{4\eta_2} - \frac{1}{2} \log(-2\eta_2) \\ h(x) &= \frac{1}{\sqrt{2\pi}} \end{aligned}$$

- Bernoulli distribution

$$f(y; \mu) = \mu^y(1 - \mu)^{1-y} = \exp \left\{ y \log \frac{\mu}{1 - \mu} + \log(1 - \mu) \right\}$$

where $\mu = \Pr(Y = 1) = E(Y)$

then

$$\begin{aligned} \theta &= \log \frac{\mu}{1 - \mu}, \phi = 1 \\ a(\phi) &= 1 \\ b(\theta) &= -\log(1 - \mu) = \log(1 + e^\theta) \\ c(y, \phi) &= c(y) = 0 \end{aligned}$$

with another formulation

$$\begin{aligned} p(x | \pi) &= \pi^x(1 - \pi)^{1-x} \\ &= \exp \left\{ \log \left(\frac{\pi}{1 - \pi} \right) x + \log(1 - \pi) \right\} \end{aligned}$$

then

$$\begin{aligned}\eta &= \frac{\pi}{1-\pi} \\ T(x) &= x \\ A(\eta) &= -\log(1-\pi) = \log(1+e^\eta) \\ h(x) &= 1\end{aligned}$$

If $Y \sim f(y; \theta, \phi)$ and ϕ is fixed, then

Property 1. $E(Y) = b'(\theta)$

Property 2. $\text{Var}(Y) = b''(\theta)a(\phi)$

Proof

$$\ell(\theta) = \log f(Y; \theta, \phi) = \frac{Y\theta - b(\theta)}{a(\phi)} + c(Y, \phi)$$

Then the score function for θ is

$$U(\theta) = \frac{\partial \ell(\theta)}{\partial \theta} = \frac{Y - b'(\theta)}{a(\phi)}$$

Since $E(U(\theta)) = 0$, we have

$$E\left(\frac{Y - b'(\theta)}{a(\phi)}\right) = 0 \implies E(Y) = b'(\theta)$$

Since $E(U(\theta)) = 0$, we have

$$\text{Var}(U(\theta)) = E\left(\frac{\partial \ell(\theta)}{\partial \theta}\right)^2 = -E\left(\frac{\partial^2 \ell(\theta)}{\partial \theta^2}\right) \implies \frac{\text{Var}(Y)}{a^2(\phi)} = \frac{b''(\theta)}{a(\phi)} \implies \text{Var}(Y) = b''(\theta)a(\phi)$$

Example: Poisson distribution

$$f(y; \mu) = \frac{e^{-\mu} \mu^y}{y!} = \exp[y \log \mu - \mu - \log y!]$$

then

$$\begin{aligned}\theta &= \log \mu \\ b(\theta) &= \mu = \exp(\theta) \\ a(\phi) &= \phi = 1\end{aligned}$$

It follows that

$$\begin{aligned}E(Y) &= b'(\theta) = \exp(\theta) = \mu \\ \text{Var}(Y) &= b''(\theta)a(\phi) = \exp(\theta) = \mu\end{aligned}$$

1.5 Likelihood theory

Likelihood definition: random variable $Y f(\tilde{Y}; \theta)$, assume n iid observations $y = (y_1, y_2, y_3, \dots, y_n)$. Then the likelihood function $L(y; \theta) = \prod f(y; \theta)$, the log-likelihood $\log L(y; \theta) = \sum f(y; \theta)$. The MLE estimator $\hat{\theta} = \text{argmax} \log L(y; \theta)$

Fisher's score function:

$$u(\theta) = \frac{\partial \log L(y; \theta)}{\partial \theta}$$

Then the MLE estimator $\hat{\theta}$ satisfies $u(\hat{\theta}) = 0$.

Fisher information: The variance of the score is defined to be the Fisher information

$$\mathcal{I}(\theta) = \text{var}(u(\theta)) = E(u(\theta)u^T(\theta)) = E\left[\left(\frac{\partial}{\partial \theta} \log f(X; \theta)\right)^2 \mid \theta\right] = -E\left[\frac{\partial^2}{\partial \theta^2} \log f(X; \theta) \mid \theta\right]$$

Asymptotic theory: the MLE estimator $\hat{\theta}$ satisfies

$$\sqrt{n}(\hat{\theta} - \theta) \rightarrow N(0, I^{-1}(\theta))$$

1.6 Likelihood for observed data

Suppose we observe a random sample of n independent observations, $O_1 = (y_1, \mathbf{x}_1^T), \dots, O_n = (y_n, \mathbf{x}_n^T)$. Our goal is to estimate the parameters θ involved in the conditional distribution of Y given \mathbf{x} , which is typically the case in regression analysis. Suppose y_i has density $f(y_i | \mathbf{x}_i, \theta)$. The joint density of \mathbf{y} is

$$f(\mathbf{y}; \theta) = \prod_{i=1}^n f_i(y_i; \theta) = \prod_{i=1}^n f(y_i | \mathbf{x}_i, \theta).$$

The likelihood (function) of θ is

$$L(\theta; \mathbf{y}) = \prod_{i=1}^n f(y_i | \mathbf{x}_i, \theta).$$

The log-likelihood of θ is

$$\ell(\theta; \mathbf{y}) = \log(L(\theta; \mathbf{y})) = \sum_{i=1}^n \log f(y_i | \mathbf{x}_i, \theta).$$

A **maximum likelihood estimator (MLE)** is a maximizer of the likelihood $L(\theta; \mathbf{y})$, denoted by $\hat{\theta}$,

$$\hat{\theta} = \operatorname{argmax}_{\theta \in \Theta} L(\theta; \mathbf{y}) = \operatorname{argmax}_{\theta \in \Theta} \ell(\theta; \mathbf{y}),$$

where Θ is the parameter space. $\hat{\theta}$ is usually obtained by solving the likelihood (score) equations. The **score function** is defined as

$$\mathbf{U}(\theta) = \frac{\partial \ell(\theta)}{\partial \theta} = \sum_{i=1}^n \frac{\partial \ell_i(\theta)}{\partial \theta} = \sum_{i=1}^n \mathbf{U}_i(\theta).$$

The **likelihood (score) equations** is defined as

$$\mathbf{U}(\theta) = 0.$$

Then $\hat{\theta}$ are zeros (solutions) of the score equations.

The **(Fisher) information matrix** is defined as

$$\begin{aligned} \mathbf{I}_n(\theta) &= -E \frac{\partial^2 \ell(\theta)}{\partial \theta \partial \theta^T} = -\sum_{i=1}^n E \frac{\partial^2 \ell_i(\theta)}{\partial \theta \partial \theta^T} = n\mathbf{I}(\theta) \\ &= E(\mathbf{U}(\theta)\mathbf{U}^T(\theta)) = \sum_{i=1}^n E(\mathbf{U}_i(\theta)\mathbf{U}_i^T(\theta)), \end{aligned}$$

$\mathbf{I}(\theta)$ is the information matrix for a *single* observation.

The **observed information matrix** is defined as

$$\begin{aligned} \mathbf{i}_n(\theta) &= -\frac{\partial^2 \ell(\theta)}{\partial \theta \partial \theta^T} = -\sum_{i=1}^n \frac{\partial^2 \ell_i(\theta)}{\partial \theta \partial \theta^T} \\ &= \mathbf{U}(\theta)\mathbf{U}^T(\theta) = \sum_{i=1}^n \mathbf{U}_i(\theta)\mathbf{U}_i^T(\theta). \end{aligned}$$

Property 1. Let θ_0 denote the true value. We have

$$E(\mathbf{U}(\theta_0)) = \mathbf{0}$$

Property 2. If $O_i = (y_i, \mathbf{x}_i^T)$ are iid, then

$$\frac{1}{n} \mathbf{i}_n(\theta) \rightarrow_p \mathbf{I}(\theta)$$

Property 3. Under some regularity conditions (e.g., the model is correctly specified), the MLE $\hat{\theta}$ has properties:

$$\text{Consistency: } \hat{\theta} \rightarrow_p \theta_0$$

$$\text{Asymptotic Normality: } \sqrt{n}(\hat{\theta} - \theta_0) \rightarrow_d \text{MVN}(0, \mathbf{I}^{-1}(\theta_0))$$