

Review for the Midterm Test

Ch 1. Generation of Random Variables

Ch 2. Optimization



Midterm Test 2021

- **Time and Date**

16:20 – 18:20,

December 20 (Monday), 120 minutes

- **Venue**

Lychee Hills Building No. 2, Room 101

- **Range**

Chapters 1 and 2

- **Assessment**

- Each assignment (5%)

- Midterm test (25%)

Key Points in the Midterm Test 2021

1. Use the inverse method to generate a discrete/continuous r.v. (20 marks);
2. Use the rejection algorithm to generate a continuous r.v. (20 marks);
3. State SIR algorithm; Use the SR, conditional sampling method to generate a multivariate r.v. (10 marks);
4. The Newton–Raphson and Fisher scoring algorithms (20 marks).
5. The EM, (ECM, MM) algorithms (30 marks)

The Policy of Closed Book Midterm Test

- Please bring one **calculator** and check the battery.
- Please bring **two** pens/pencils in case one is not available.
- You can prepare anything on **one side** of an A4 paper and bring it with you to the test venue.
- You are not allowed to bring any other material (including **iPhone/iPad**) to the test venue.

Ch 1. Generation of Random Variables

- §1.1 The Inversion Method
- §1.2 The Grid Method
- §1.3 The Rejection Method
- §1.4 The SIR Method
- §1.5 The SR Method
- §1.6 The CS Method



§1.1 Three Key Points to The Inversion Method

1. The algorithm for the continuous r.v.:

Step 1. Draw $U = u$ from $U(0, 1)$;

Step 2. Return $x = F^{-1}(u)$.

2. How to generate samples from a discrete distribution using the inversion method (please review all examples in Subsection 1.1.2).

3. The built-in R function:

`sample(x, N, prob = p, replace = T/F)`
(please review Appendix A.1.1).

Remarks on the Inversion Method

- 1° **Applicability:** If F^{-1} has an explicit expression, the inversion method is the most efficient method to generate a random sample of X from $F(\cdot)$.
- 2° **Inapplicability:** If F^{-1} is not available analytically, the inversion method may not be efficient.
- 3° The inversion method is a special case of the **stochastic representation** (SR) method:

$$X \stackrel{d}{=} F^{-1}(U) \stackrel{d}{=} F^{-1}(1 - U).$$

§1.2 The Grid Method

- The **essence** of the grid method is to generate a finite discrete distribution:

$$X \sim \mathbf{FDiscrete}_d(\{x_i\}, \{p_i\}).$$

where $\{x_i\}_{i=1}^d$ is a set of grid points, covering the support \mathcal{S}_X ,

- and

$$p_i = \frac{f_X(x_i)}{\sum_{j=1}^d f_X(x_j)}, \quad i = 1, \dots, d.$$

Remarks on the Grid Method

- 1° **Applicability:** If (i) the normalizing constant of $f_X(x)$ is either known or unknown; and (ii) the support of X is a finite interval, i.e., $\mathcal{S}_X = [a, b]$, where $-\infty < a$ and $b < +\infty$, then the grid method is an efficient method to generate a random sample of X from $f_X(\cdot)$.
- 2° **Inapplicability:** If \mathcal{S}_X is an infinite interval, the grid method cannot be applied.

§1.3 Four Key Points to The Rejection Method

1. **Several basic notions:** Target density $f(x)$; Envelope constant c ; Envelope density $g(x)$; Acceptance probability $1/c$.
2. **The algorithm:**

Step 1. Draw $U \sim U(0, 1)$ and independently draw $Y \sim g(\cdot)$;

Step 2. If $U \leq \frac{f(Y)}{cg(Y)}$, return $X = Y$; otherwise, go to Step 1.

§1.3 Four Key Points to The Rejection Method (Cont'd)

3. Theoretical justification (see page 17)
4. How to find the optimal c_{opt} among

$$c_{\theta} = \max_{x \in \mathcal{S}_X} \frac{f(x)}{g_{\theta}(x)}, \quad (1.13)$$

where $\{g_{\theta}(x) : \theta \in \Theta\}$ is a family of the candidate envelope densities indexed by a parameter θ .

(Please review all examples in Subsection 1.3.3)

§1.4 Two Points to the SIR Method

1. **Two basic notions:** Target density $f(x)$; Important sampling density $g(x)$.
2. **The algorithm:**

Step 1. Generate $X^{(1)}, \dots, X^{(J)} \stackrel{\text{iid}}{\sim} g(\cdot)$;

Step 2. Select a subset $\{X^{(k_i)}\}_{i=1}^m$ from $\{X^{(j)}\}_{j=1}^J$ via resampling **without replacement** from the discrete distribution on $\{X^{(j)}\}$ with probabilities $\{\omega_j\}$.

§1.5 A Key Point to The SR Method

1. How to prove that a given multivariate random vector can be represented by other independent r.v.s.

- Prove that $\mathbf{x} = (X_1, \dots, X_d)^\top \sim \text{Dirichlet}(a_1, \dots, a_d)$ has the following SR

$$X_i \stackrel{d}{=} (1 - W_{i-1}) \prod_{j=i}^{d-1} W_j, \quad i = 1, \dots, d-1,$$

$$X_d \stackrel{d}{=} 1 - W_{d-1}$$

where $W_0 \equiv 0$, $\{W_i\}_{i=1}^{d-1} \stackrel{\text{ind}}{\sim} \text{Beta}(a_1 + \dots + a_i, a_{i+1})$.

- You only need to prove that

$$f(\mathbf{x}_{-d}) = f(x_1, \dots, x_{d-1}) = \frac{\Gamma(\sum_{i=1}^d a_i)}{\prod_{i=1}^d \Gamma(a_i)} \prod_{i=1}^d x_i^{a_i-1}.$$

- Note that

$$f(\mathbf{x}_{-d}) = \left[\prod_{i=1}^{d-1} g_i(w_i) \right] \cdot |J(\mathbf{w}_{-d} \rightarrow \mathbf{x}_{-d})|,$$

where the Jacobian determinant is

$$J(\mathbf{w}_{-d} \rightarrow \mathbf{x}_{-d}) = \frac{\partial(w_1, \dots, w_{d-1})}{\partial(x_1, \dots, x_{d-1})}.$$

§1.6 The Conditional Sampling Method

- Let $\mathbf{x} = (X_1, \dots, X_d)^\top$ and its density $f_{\mathbf{x}}(\mathbf{x})$ can be factorized as

$$f_{\mathbf{x}}(\mathbf{x}) = f_1(x_1) \left\{ \prod_{i=2}^d f_i(x_i | x_1, x_2, \dots, x_{i-1}) \right\}. \quad (1.33)$$

- To generate \mathbf{x} from $f_{\mathbf{x}}(\mathbf{x})$, we only need to generate x_1 from the marginal density $f_1(x_1)$, then to generate x_i sequentially from the conditional density $f_i(x_i | x_1, x_2, \dots, x_{i-1})$. (Review all examples in Section 1.6)

Ch 2. Optimization

- §2.1 Rate of Convergence
- §2.2 The NR Algorithm
- §2.3 The EM Algorithm
- §2.4 The ECM Algorithm
- §2.5 MM Algorithms



§2.1 Rate of convergence

1. Let an **EM/MM algorithm** can be represented by $\theta^{(t+1)} = h(\theta^{(t)})$, then, the rate of convergence of the EM/MM algorithm is defined by (2.9), i.e.,

$$c = \lim_{t \rightarrow \infty} \frac{|\theta^{(t+1)} - \hat{\theta}|}{|\theta^{(t)} - \hat{\theta}|} = |h'(\hat{\theta})|,$$

where $\hat{\theta}$ is the MLE of θ .

2. Calculate the value of the rate of convergence of the EM algorithm in Example 2.6

3. Let an **NR/FS algorithm** can be represented by $\theta^{(t+1)} = f(\theta^{(t)})$, then, the rate of convergence of the NR/FS algorithm is defined by

$$c = \lim_{t \rightarrow \infty} \frac{|\theta^{(t+1)} - \hat{\theta}|}{|\theta^{(t)} - \hat{\theta}|^2} = \frac{1}{2} |f''(\hat{\theta})|,$$

where $\hat{\theta}$ is the MLE of θ .

4. Calculate the value of the rate of convergence of the NR/FS algorithm in Example 2.6

§2.2 The NR Algorithm

1. **Several important notions:** Score vector $\nabla\ell(\boldsymbol{\theta}|Y_{\text{obs}})$; Observed information matrix $\mathbf{I}(\boldsymbol{\theta}|Y_{\text{obs}})$; Fisher/expected information matrix $\mathbf{J}(\boldsymbol{\theta})$; Estimated asymptotic covariance matrix $\widehat{\text{Cov}}(\hat{\boldsymbol{\theta}})$.

2. **The NR algorithm**

$$\boldsymbol{\theta}^{(t+1)} = \boldsymbol{\theta}^{(t)} + \mathbf{I}^{-1}(\boldsymbol{\theta}^{(t)}|Y_{\text{obs}})\nabla\ell(\boldsymbol{\theta}^{(t)}|Y_{\text{obs}}). \quad (2.13)$$

3. **Application to logistic regression.**
(Please review §2.2.5)

4. How to apply the NR/FS algorithms to the Poisson regression. [See Ex. T4.3]

Let $Y_{\text{obs}} = \{y_i\}_{i=1}^n$ and consider the following Poisson regression

$$Y_i \stackrel{\text{ind}}{\sim} \text{Poisson}(\lambda_i),$$
$$\log(\lambda_i) = \mathbf{x}_{(i)}^\top \boldsymbol{\theta}, \quad 1 \leq i \leq n,$$

then the log-likelihood function for $\boldsymbol{\theta}$ is

$$\ell(\boldsymbol{\theta} | Y_{\text{obs}}) = c + \sum_{i=1}^n \left\{ y_i (\mathbf{x}_{(i)}^\top \boldsymbol{\theta}) - \exp(\mathbf{x}_{(i)}^\top \boldsymbol{\theta}) \right\}.$$

§2.3 The EM Algorithm

1. Summary of the EM algorithm

- Augment the observed data Y_{obs} with latent variables Z .
- **M-Step:** Find the complete-data log-likelihood function $\ell(\boldsymbol{\theta}|Y_{\text{obs}}, z)$ and derive the complete-data MLE $\hat{\boldsymbol{\theta}}$.
- **E-Step:** Find the conditional predictive distribution $f(z|Y_{\text{obs}}, \boldsymbol{\theta})$ and compute the conditional expectation $E(Z_i|Y_{\text{obs}}, \boldsymbol{\theta})$ or/and $E(Z_i^2|Y_{\text{obs}}, \boldsymbol{\theta})$.

2. The ascent property of the EM

- The definition of the Q function

$$\begin{aligned} Q(\boldsymbol{\theta}|\boldsymbol{\theta}^{(t)}) &= E\{\ell(\boldsymbol{\theta}|Y_{\text{obs}}, Z)|Y_{\text{obs}}, \boldsymbol{\theta}^{(t)}\} \\ &= \int_{\mathbb{Z}} \ell(\boldsymbol{\theta}|Y_{\text{obs}}, z) \times f(z|Y_{\text{obs}}, \boldsymbol{\theta}^{(t)}) \, dz, \end{aligned}$$

- The original definition of the EM

$$\boldsymbol{\theta}^{(t+1)} = \arg \max_{\boldsymbol{\theta} \in \Theta} Q(\boldsymbol{\theta}|\boldsymbol{\theta}^{(t)}). \quad (2.18)$$

- Using the Kullback–Leibler divergence to prove that

$$\ell(\boldsymbol{\theta}^{(t+1)}|Y_{\text{obs}}) - \ell(\boldsymbol{\theta}^{(t)}|Y_{\text{obs}}) \geq 0.$$

3. Missing information principle

- The observed information, the complete information, and the missing information:

$$\mathbf{I}_{\text{obs}} = \mathbf{I}_{\text{com}} - \mathbf{I}_{\text{mis}}. \quad (2.38)$$

§2.5 MM Algorithms

1. The MM idea

- **Minorization.** The function $Q(\boldsymbol{\theta}|\boldsymbol{\theta}^{(t)})$ is said to minorize $\ell(\boldsymbol{\theta}|Y_{\text{obs}})$ at $\boldsymbol{\theta} = \boldsymbol{\theta}^{(t)}$ if

$$Q(\boldsymbol{\theta}|\boldsymbol{\theta}^{(t)}) \leq \ell(\boldsymbol{\theta}|Y_{\text{obs}}) \quad \forall \boldsymbol{\theta}, \boldsymbol{\theta}^{(t)} \in \Theta, \quad (2.43)$$

$$Q(\boldsymbol{\theta}^{(t)}|\boldsymbol{\theta}^{(t)}) = \ell(\boldsymbol{\theta}^{(t)}|Y_{\text{obs}}). \quad (2.44)$$

- **Definition of an MM algorithm.**

$$\boldsymbol{\theta}^{(t+1)} = \arg \max_{\boldsymbol{\theta} \in \Theta} Q(\boldsymbol{\theta}|\boldsymbol{\theta}^{(t)}), \quad (2.45)$$

- **The ascent property of an MM.**

2. The QLB algorithm

- **A key condition.** There exists a positive definite matrix $\mathbf{B} > 0$ such that

$$\nabla^2 \ell(\boldsymbol{\theta} | Y_{\text{obs}}) + \mathbf{B} \geq 0 \quad \forall \boldsymbol{\theta} \in \Theta. \quad (2.47)$$

- **Definition of the Q function.**

$$Q(\boldsymbol{\theta} | \boldsymbol{\theta}^{(t)}) = \ell(\boldsymbol{\theta}^{(t)} | Y_{\text{obs}}) + (\boldsymbol{\theta} - \boldsymbol{\theta}^{(t)})^\top \nabla \ell(\boldsymbol{\theta}^{(t)} | Y_{\text{obs}}) - \frac{1}{2} (\boldsymbol{\theta} - \boldsymbol{\theta}^{(t)})^\top \mathbf{B} (\boldsymbol{\theta} - \boldsymbol{\theta}^{(t)}), \quad \boldsymbol{\theta} \in \Theta. \quad (2.48)$$

- **The Definition of the QLB algorithm.**

$$\boldsymbol{\theta}^{(t+1)} = \boldsymbol{\theta}^{(t)} + \mathbf{B}^{-1} \nabla \ell(\boldsymbol{\theta}^{(t)} | Y_{\text{obs}}). \quad (2.49)$$

3. De Pierro's algorithm

- **A key condition.** The log-likelihood function be of the form

$$\ell(\boldsymbol{\theta}|Y_{\text{obs}}) = \sum_{i=1}^m f_i(\mathbf{x}_{(i)}^\top \boldsymbol{\theta}),$$

where $\{f_i\}_{i=1}^m$ are twice continuously differentiable and strictly concave functions defined in \mathbb{R} .

- **Definition of the Q function.**

$$Q(\boldsymbol{\theta}|\boldsymbol{\theta}^{(t)}) = \sum_{i=1}^m \sum_{j \in \mathbb{J}_i} \lambda_{ij} f_i(\lambda_{ij}^{-1} x_{ij}(\theta_j - \theta_j^{(t)}) + \mathbf{x}_{(i)}^\top \boldsymbol{\theta}^{(t)}), \quad (2.52)$$

4. Jensen's Inequality

- Let $\varphi(\cdot)$ be concave. If X is a r.v. taking values in the domain of $\varphi(\cdot)$, then

$$\varphi[E(X)] \geq E[\varphi(X)],$$

provided that $E(X)$ and $E[\varphi(X)]$ exist.

- **Discrete version.** For any concave function $f(\cdot)$,

$$f\left(\sum_{i=1}^n \alpha_i x_i\right) \geq \sum_{i=1}^n \alpha_i f(x_i),$$

where $\alpha_i \geq 0$ and $\sum_{i=1}^n \alpha_i = 1$.

5. The supporting hyperplane inequality

- If $g(x)$ is convex; i.e., $g''(x) \geq 0$, we have

$$g(x) \geq g(x_0) + (x - x_0)g'(x_0).$$

- Please review all questions in Assignments 1-3.
- Please review all questions in Tutorials 1-7.

End of the Review



易游人
YOUYOU

Pkyangxing 51.com