

STATISTICS 743 (Part 2)

Foundations of Statistics

Angelo J. Canty

Course Information

Lectures Monday & Thursday 12:30-2:00 in BSB-B140.

Office : Hamilton Hall 209

Phone : (905) 525-9140 extn 27079

E-mail : cantya@mcmaster.ca

Website : All lecture notes and announcements about the course will be made on the website

<http://www.math.mcmaster.ca/canty/teaching/stat743>

Course Information

Office Hours I will generally be in my office for 30–45 minutes after each class. Appointments can also be made by e-mail. For short questions, you can drop by my office and if I am there and free I will answer your question.

Assignments 4 assignments will be worth 60% of your mark (15% each).

Final Exam 3 hour written exam in April will be worth the other 40% of your mark.

Course Grade Based on an equally weighted average of marks in both terms.

Topics to be covered

Week 1 Simulation and Monte Carlo Methods

Week 2 Point Estimation

Week 3 Computation of the MLE

Weeks 4–5 Hypothesis testing

Weeks 6–7 Interval estimation

Weeks 8–9 Bayesian inference

Weeks 10–11 Consistency, efficiency and Asymptotics.

Weeks 12–13 Bootstrap and related techniques

Monte Carlo Methods

- * It is quite common in statistics that we need to evaluate complex integrals.
- * Such integrals are often quite difficult (or indeed impossible) to evaluate in closed form.
- * Monte Carlo methods are a way of approximating the value of an integral using large samples of random variables.
- * These samples of random variables are typically computer generated.

Monte Carlo Integration

- * Suppose that we need to evaluate

$$I = \int_{\mathcal{A}} h(x) dx$$

- * Let f be a probability density function with support \mathcal{A} .

- * Then we can write

$$I = \int_{\mathcal{A}} \frac{h(x)}{f(x)} f(x) dx = \int_{\mathcal{A}} g(x) f(x) dx$$

- * Now if Y is a random variable with pdf f then we have

$$I = \mathbb{E}[g(Y)].$$

Monte Carlo Integration

- * Hence if we have $Y_1, \dots, Y_N \stackrel{iid}{\sim} f$, the Weak Law of Large Numbers tells us that

$$\hat{I} = \frac{1}{N} \sum_{i=1}^N g(Y_i) \xrightarrow{p} I.$$

- * We can therefore use \hat{I} to approximate I very well for large enough N .
- * Furthermore we can use estimate the variability in \hat{I} using the sample variance of the random sample $g(Y_1), \dots, g(Y_N)$ divided by N .
- * Since N is totally in our control we can choose N to be large enough to make the variability as low as we desire.

Generating Uniform Random Variates

- * Computers are unable to generate random numbers.
- * They can, however, be used to generate **pseudo-random** numbers.
- * These are sequences of numbers which are generated from a deterministic algorithm but which behave like a sequence of *iid* random variates.

Generating Uniform Random Variates

- * Generally a computer uses an algorithm to generate a sequence of integers x_0, x_1, x_2, \dots such that

$$x_t \in \{0, 1, \dots, M - 1\}$$

- * The generated values are then rescaled to be approximately Uniform(0, 1)

$$U_t = \frac{x_t}{M}$$

- * Many algorithms have been proposed. Most are of the form

$$x_t = g(x_{t-1})$$

- * The function g is chosen so that the resulting sequence resembles a random sequence.

Common Random Number Generators

Definition 1.1

A *linear congruential generator* on $\{0, 1, \dots, M - 1\}$ is a sequence of integers defined by

$$x_{t+1} = (ax_t + b) \bmod M$$

Definition 1.2

For a given $k \times k$ matrix T whose entries are all either 0 or 1, the associated *shift register generator* is

$$x_{t+1} = Tx_t \bmod 2$$

where x_t and x_{t+1} are binary representations of the corresponding numbers.

Combination Generators

- * Most modern generators use two or more parallel generators and return a sum or product of the results modulo the largest integer representable.
- * One such generator is George Marsaglia's KISS (Keep It Simple, Stupid) generator which uses one linear congruential generator and two shift register generators, returning their sum modulo 2^{32} .
- * The current state of the art generator is generally believed to be the Mersenne-Twister Generator and this is the default in R.
- * No test of uniformity or randomness has yet been found that the sequence from this generator does not pass!

Reproducibility of Generated Random Numbers

- * An important concept is that of a **Random Seed**.
- * The seed specifies initial conditions for the algorithm which generates the random variates.
- * Since the algorithm is actually deterministic, setting the same random seed will result in the same sequence of random variates.
- * When using simulation to compare two or more methods, we will often want to use the same random numbers for each method to remove the simulation variability between methods and so give a better comparison.

Generating Non-Uniform Random Numbers

- * Uniform random variates are rarely what we need for simulation or Monte Carlo inference.
- * They are, however, the building blocks for generating random variates from any other distributions.
- * Much of this is based on the following theorem

Theorem 1.1 (Probability Integral Transform)

Suppose that $U \sim \text{Uniform}(0, 1)$ and that F is a continuous cdf with unique inverse F^{-1} . Then the random variable

$$Y = F^{-1}(U)$$

has a distribution with cdf F .

Generating Discrete Random Variates

- * The method as described above requires that we have a continuous cdf.
- * A similar technique can also be used to generate discrete random variables.
- * Suppose that $p(y)$ is the probability mass function and the support of the random variable is $\mathcal{Y} = \{y : p(y) > 0\}$. Then we can define the inverse cdf as

$$F^{-1}(u) = \min\{y \in \mathcal{Y} : F(y) \geq u\}$$

- * Then if $U \sim \text{uniform}(0, 1)$, the random variable $Y = F^{-1}(u)$ will be distributed with probability mass function $p(y)$.

Special Methods

- * In many cases, the inverse cdf is not available in closed form and so this method cannot be used. We can often, however, use algorithms based on transformations for such situations.

- * Suppose that U_1 and U_2 are two independent Uniform(0,1) random variables then it is easy to show that

$$Y_1 = \sqrt{-2 \log U_1} \sin(2\pi U_2) \quad \text{and} \quad Y_2 = \sqrt{-2 \log U_1} \cos(2\pi U_2)$$

are independent standard normal random variables.

- * This is known as the **Box-Muller Algorithm**.

Accept/Reject Algorithm

- * A more general technique which is useful when the inverse cdf method cannot be applied is called the **Accept/Reject Algorithm**.
- * This method relies on generating a different random variable V which has the same support as the required variable Y .
- * We also require that the ratio of densities is bounded by a known constant

$$M = \sup_y \frac{f_Y(y)}{f_V(y)} < \infty$$

Accept/Reject Algorithm

1. Calculate $M = \sup_y f_Y(y)/f_V(y)$.
2. Generate $V \sim f_V$ and independently $U \sim \text{Uniform}(0, 1)$.
3. If

$$U < \frac{f_Y(V)}{M f_V(V)}$$

then set $Y = V$. Otherwise discard U and V and return to step 2.

Markov Chain Monte Carlo Methods

- * Many of the methods described so far are not very useful for generating multivariate random variates.
- * Markov Chain Monte Carlo methods are now widely used in these settings.
- * The methods work on the idea of constructing a Markov chain which has a stationary distribution equal to the distribution of interest.
- * Under certain conditions, the distribution of the elements in such a chain will converge to this stationary distribution.

Markov Chain Monte Carlo Methods

- * These algorithms start with some initial value for the random variable of interest.
- * They then run a carefully constructed Markov chain starting from that initial value for a sufficiently long time.
- * It is not always easy to know how long the chains should be run but various diagnostics have been proposed.
- * Any observations in the chain after this **burn-in** period may be considered as (at least approximately) distributed with the stationary distribution.

Metropolis–Hastings Algorithm

- * First introduced in statistical physics in 1954 by Metropolis et al. Statistical properties shown by Hastings in 1970.
- * It is basically a Markov chain version of the accept/reject algorithm.
- * Random variates are generated from some candidate distribution conditional on the current state of the chain and then either the new state is accepted or rejected in which case the chain stays where it is.

Metropolis–Hastings Algorithm

Suppose we wish to sample $Y \sim f_Y$.

First initialize the chain with some value $Y^{(0)}$.

Then for $t = 1, 2, \dots$ we generate $Y^{(t)}$ by

1. Generate $V^{(t)} \sim f_{V|Y}(v | Y^{(t-1)})$.
2. Calculate the acceptance probability

$$\rho_t = \min \left\{ \frac{f_Y(V^{(t)})}{f_Y(Y^{(t-1)})} \times \frac{f_{V|Y}(Y^{(t-1)} | V^{(t)})}{f_{V|Y}(V^{(t)} | Y^{(t-1)})}, 1 \right\}$$

3. Generate $U_t \sim \text{Uniform}(0, 1)$ and set

$$Y^{(t)} = \begin{cases} V^{(t)} & \text{if } U_t \leq \rho_t \\ Y^{(t-1)} & \text{if } U_t > \rho_t \end{cases}$$

Independence Metropolis–Hastings Algorithm

- * It is often convenient to generate $V^{(t)}$ from the same distribution at every iteration.
- * In this case we have $f_{V|Y}(v \mid Y^{(t-1)}) = f_V(v)$ and so the acceptance probability becomes

$$\begin{aligned}\rho_t &= \min \left\{ \frac{f_Y(V^{(t)})}{f_Y(Y^{(t-1)})} \times \frac{f_V(Y^{(t-1)})}{f_V(V^{(t)})}, 1 \right\} \\ &= \min \left\{ \frac{f_Y(V^{(t)})}{f_V(V^{(t)})} \times \frac{f_V(Y^{(t-1)})}{f_Y(Y^{(t-1)})}, 1 \right\}\end{aligned}$$

Random Walk Metropolis–Hastings Algorithm

- * Another special case is where $f_{V|Y}(v | y) = f_Z(v - y)$ where f_Z is a distribution symmetric about 0.
- * We generate $Z^{(t)} \sim f_Z$ and set $V^{(t)} = Y^{(t-1)} + Z^{(t)}$.
- * In stochastic processes this is called a random walk.
- * The acceptance probability for the Metropolis–Hastings algorithm then becomes

$$\rho_t = \min \left\{ \frac{f_Y(V^{(t)})}{f_Y(Y^{(t-1)})}, 1 \right\}$$

Gibbs Sampler

- * The **Gibbs Sampler** (Geman & Geman, 1984) is designed to generate observations from a complex multivariate distribution.
- * The Markov chain is constructed by considering the univariate conditional distributions.
- * Suppose that the random vector of interest is $\mathbf{Y} = (Y_1, \dots, Y_d)$ and that we can generate observations from the **full conditional distributions**

$$f_j(y \mid Y_{-j} = y_{-j}) = f_{y_j|y_{-j}}(y \mid Y_{-j} = y_{-j}) \quad j = 1, \dots, d$$

where $Y_{-j} = (Y_1, \dots, Y_{j-1}, Y_{j+1}, \dots, Y_d)$.

Gibbs Sampler

Initialise the chain to some value $\mathbf{Y}^{(0)} = (Y_1^{(0)}, \dots, Y_d^{(0)})$.

For $t = 1, 2, \dots$

- 1 Generate $Y_1^{(t)}$ from $f_1(y_1 | Y_2^{(t-1)}, \dots, Y_d^{(t-1)})$.
- 2 Generate $Y_2^{(t)}$ from $f_2(y_2 | Y_1^{(t)}, Y_3^{(t-1)}, \dots, Y_d^{(t-1)})$.
- \vdots
- j Generate $Y_j^{(t)}$ from $f_j(y_j | Y_1^{(t)}, \dots, Y_{j-1}^{(t)}, Y_{j+1}^{(t-1)}, \dots, Y_d^{(t-1)})$.
- \vdots
- d Generate $Y_d^{(t)}$ from $f_d(y_d | Y_1^{(t)}, \dots, Y_{d-1}^{(t)})$.

Then we set $\mathbf{Y}^{(t)} = (Y_1^{(t)}, \dots, Y_d^{(t)})$.