

Monte Carlo Inference

Part III Course, Lent 2010

Revision Notes

Daniel Guetta

guetta@cantab.net

Generally Useful R Code

Distributions

- Normal
 - `rnorm(n, m, s)` generates n data points from $N(m, s^2)$
 - `pnorm(x, m, s) = $\mathbb{P}(N(m, s^2) \leq x)$`
 - `dnorm(x, m, s) = $f_X(x)$` if $X \sim N(m, s^2)$. If the argument `log = TRUE` is specified, `log($f_X(x)$)` is returned.
- Gamma
 - `rgamma(n, a, rate=b)` generates n data points from $\Gamma(a, b)$
 - Others as above
- Uniform
 - `runif(n, a, b)` generates n data points from $U(a, b)$
 - Others as above
- Students' t
 - `rt(n, df=1)` generates n data points from a t distribution with 1 degree of freedom (ie: a Cauchy distribution).
 - Others as above
- Chi-squared
 - `rchisq(n, df=5)` generates n data points from a chi-squared distribution with 5 degrees of freedom.
 - Others as above

Graphics

- `par()` is used to modify settings in the current graphics environment. Possible settings:
 - `mfrow = c(rows, columns)` sets the numbers of rows and columns of graphs
 - `cex.main=s`, `cex.lab=s`, `cex.axis=s` respectively set the title, axes labels and axes font sizes to s
- `plot(x, y)` creates a new plot area and plots two vectors against each other. Possible settings
 - `type = "l"` plots a line (by default, individual points are plotted)

- `lwd = 2` makes the line (or points) thicker
 - `col = "blue"` makes the line blue
 - `ylim = c(miny, maxy)` sets the limits on the axes [could also use `range` function below]
 - `lty = 2` or `3` give different kinds of dashing.
 - `main = "title"` sets the title to "main"
- `lines(x, y)` adds a line to an existing plot. Possible settings as for `plot`.
- `points(x, y)` adds points to an existing plot. Possible settings as for `plot`.
- `abline(v=1)` adds a vertical line at $x = 1$, and `abline(h=1)` adds a horizontal line at $y = 1$. Possible settings as for `plot`.
- `hist(data)` plots a histogram of the input data. Possible settings as for `plot`, as well as
 - `freq = FALSE` normalises the area of the histogram to 1. The default, `freq = TRUE`, simply plots counts.
- `text(x, y, "text")` adds text to a plot. `col` and `cex` are possible options for the colour and font size. Note that `y = 0` puts the text straight on the x -axis.
- `d <- dev.cur` sets a handle for the current graphics window. A new graphics window can then be created using `x11()`, and the old graphics window can be returned to using `dev.set(d)`.

Data handling

- `plot(density(x))` plots the empirical density of x . Very useful if trying to pictorially show the efficiency of various estimators; simulate the value of interest a number of times, and plot the density each time to see what it looks like.
- `data.frame(a=aData, b=bData)` creates a data frame with variable `a` taking values in `aData`, etc... The data can consist of vectors; R will add suffices accordingly.
- `summary(dataFrame)` gives all kinds of useful statistics on the items in the data frame.

- `boxplot(dataFrame)` plots a box-and-whisker plot showing the spread of data in each of the variables in the data frame.
- `quantiles(x, v)` finds the v quantile of the data x . For example, if $v = 0.5$, it finds the median. Extremely useful in approximating confidence intervals.

Vectors, matrices, etc...

- Vector numbering starts from 1, not 0.
- `x <- seq(a,b,length=c)` fills x with c values between a and b . Another possible argument is `by=c`, which increments by c each time.
- `rep(a,b)` produces a vector containing b instances of a . a can be set to `NA` to create an empty array.
- `range(vector)` gives a column containing the lowest and highest value of a vector. Useful when setting the limits of an axis.
- `rev(x)` reverses the input vector.
- `length(x)` gives the number of items in a vector.
- `x>2` returns a vector of the same length as x containing `TRUE` wherever the condition is met and `FALSE` otherwise. Can be used in a number of ways
 - `x[x > 2]` returns a vector containing all items in x greater than 2.
 - `mean[x > 2]` assigns 1 to each component of x greater than 2, and 0 to others, and finds the mean of these numbers.
- `apply(matrix, index, <operation>)` applies the operation `<operation>` to rows (if `index = 1`) or columns (if `index = 2`) of `matrix`.
- `matrix(data, nrow, ncol, byrow = FALSE)` creates a matrix and fills it with `data`, column by column (unless `byrow = TRUE`).
- `t(m)` transposes the matrix m .
- `x %*% y` performs the matrix multiplication xy .
- `prod(x, na.rm=FALSE)` finds the product of the elements in the vector x . If `na.rm = TRUE`, missing values are removed.
- `Numeric(0)` creates an empty vector.

Control blocks

- `for (i in 1:5) { }` loops the item in the braces 5 times.

Other

- `x <- readline("prompt")` prompts the user and inserts their response as a string into `x`.
- `paste(string1,string2,sep=".")` concatenates the strings in the argument, adding `.` between each string.
- `help(functionName)` gives help.

Random Number Generation

Assume we have an infinite supply of random numbers U that are distributed randomly over $[0,1]$. We discuss general methods for using these numbers to generate random numbers from other, more sophisticated distributions.

Method of Inversion

If X has a continuous CDF F , then $U = F(X) \sim U(0,1)$ ¹, and so $X = F^{-1}(U)$.

This is the basis of the method of inversion:

Method of inversion: To generate a sample x from a distribution with CDF F

1. Simulate $u \sim U(0,1)$
2. Set $x = F^{-1}(u)$

Proof: Consider that

$$\mathbb{P}(X \leq x) = \mathbb{P}(F^{-1}(U) \leq x)$$

since F is a strictly increasing², monotonic and continuous function of x , we can write

$$= \mathbb{P}(U \leq F(x))$$

since $U \sim U(0, 1)$, $P(U \leq ?) = ?$, and so

$$= F(x)$$

as desired. ■

Notes:

- The inverse CDF relationship exists between any two continuous random variables; $X = F_X^{-1}(F_Y(Y))$. This method is therefore often used with $Y =$ normal distribution.
- Even when F^{-1} exists in closed form, it may be more computationally intensive to calculate than some alternative methods for generating random numbers.

¹ To see why, consider that $\mathbb{P}(F(P) \leq p) = \mathbb{P}(P \leq F^{-1}(p)) = F(F^{-1}(p)) = p$.

² For functions that are not *strictly* increasing, use a generalised inverse.

- When F^{-1} does not exist in closed form, the method of inversion can still be used by solving $F_X(x) - u = 0$ numerically.

This method can also be used for discrete distributions

Method of inversion (discrete distributions):

Consider a distribution with levels m_j and $\mathbb{P}(m_j) = p_j$. The CDF is

$$F_j = \sum_{k=1}^j p_k$$

so that $p_j = F_j - F_{j-1}$. To generate a sample x from this distribution

1. Simulate $u \sim U(0,1)$
2. Set $x = m_j$ if $F_{j-1} < u < F_j$

Proof: Consider that

$$\mathbb{P}(X = j) = \mathbb{P}(F_{j-1} < U \leq F_j)$$

since $U \sim U(0, 1)$, $P(a < U \leq b) = b - a$, and so

$$\begin{aligned} &= P_j - P_{j-1} \\ &= p_j \end{aligned}$$

as desired. ■

In practice, the algorithm above can be applied by first simulating u and then calculating each $\mathbb{P}(m_j)$ for each j , until the sum of all probabilities calculated is greater or equal to u .

This method can also be used for mixture densities

Method of inversion (mixture distributions) Consider a distribution with mixture density

$$f = \sum_{i=1}^k w_i f_i$$

where the f_i are PDFs, $w_i > 0$ and $\sum w_i = 1$. To sample from this distribution:

1. Choose $I = i$ with probability w_i
2. Sample from f_I

Rejection Sampling

Rejection sampling: Suppose it is difficult to sample directly from a density f , but that we have a **majorising** or **envelope** density g from which it is easy to sample and a constant $M \in [1, \infty]$ such that

$$f(x) \leq Mg(x) \quad \forall x \in \mathbb{R}$$

Or alternatively

$$\sup_{x \in \mathbb{R}} \left(\frac{f(x)}{g(x)} \right) < \infty$$

To generate a sample x from the distribution with density g :

1. Generate $y \sim g$ and independent $u \sim U(0,1)$
2. If $u > \frac{f(y)}{Mg(y)}$, return to step 1.
3. Return $x = y$.

Proof: We have

$$\begin{aligned} \mathbb{P}(X \leq x) &= \mathbb{P}(Y \leq x \mid Y \text{ is accepted}) \\ &= \frac{\mathbb{P}(Y \leq x \text{ and } Y \text{ is accepted})}{\mathbb{P}(Y \text{ is accepted})} \end{aligned}$$

But we know that

$$\begin{aligned} \mathbb{P}(Y \text{ is accepted}) &= \mathbb{P}\left(U \leq \frac{f(Y)}{Mg(Y)}\right) \\ &= \int_{-\infty}^{\infty} \mathbb{P}\left(U \leq \frac{f(Y)}{Mg(Y)} \mid Y = y\right) \mathbb{P}(Y = y) \, dy \\ &= \int_{-\infty}^{\infty} \mathbb{P}\left(U \leq \frac{f(y)}{Mg(y)}\right) g(y) \, dy \\ &= \int_{-\infty}^{\infty} \int_0^{\frac{f(y)}{Mg(y)}} g(y) \, du \, dy \\ &= \int_{-\infty}^{\infty} \frac{f(y)}{Mg(y)} g(y) \, dy \\ &= \frac{1}{M} \int_{-\infty}^{\infty} f(y) \, dy \end{aligned}$$

Similarly

$$\begin{aligned} \mathbb{P}\left(\begin{array}{l} Y \leq x \text{ and} \\ Y \text{ is accepted} \end{array}\right) &= \int_{-\infty}^x \mathbb{P}\left(U \leq \frac{f(y)}{Mg(y)}\right) g(y) \, dy \\ &= \frac{1}{M} \int_{-\infty}^x f(y) \, dy \end{aligned}$$

as desired. ■

Note that the denominator in the final expression normalises the function, so we only need to know f up to a multiplicative constant.

Note that the proportion of “accepted” trials is given by

$$\pi = \frac{\text{Area under } f}{\text{Area under } g} = \frac{1}{M}$$

So we want to choose M as small as possible subject to $M \geq f(x)/g(x)$. Since M always needs to be bigger than the RHS, the smallest value M can take is the supremum of the RHS. So the optimal M is given by

$$M^* = \sup_{x \in \mathbb{R}} \left(\frac{f(x)}{g(x)} \right)$$

Sometimes, we may choose a *family* of enveloping functions characterised by a parameter β , say. In that case, minimizing M^* with respect to β provides the best choice of function.

EXAMPLE: Let $\alpha \in (1, \infty)$, and consider sampling from $\Gamma(\alpha, 1)$

$$f(x) = \frac{x^{\alpha-1} e^{-x}}{\Gamma(\alpha)} \quad x \in (0, \infty)$$

We may be tempted to choose a majoring function $g(x) = x^{\alpha-1}$, but that’s no good, because it is not bounded, and therefore is not a distribution. Instead, we choose the family

$$g_\beta(x) = \beta e^{-\beta x}$$

and we note that

$$\sup_{x \in \mathbb{R}^+} \left(\frac{f(x)}{g_\beta(x)} \right) = \frac{\left(\frac{\alpha-1}{1-\beta}\right)^{\alpha-1} e^{-(\alpha-1)}}{\beta \Gamma(\alpha)} < \infty$$

So g does indeed work as a majoring function. We also note that the minimum of this supremum is attained at $\beta^* = \frac{1}{\alpha}$.

EXAMPLE: Consider sampling from the distribution

$$f(x) = e^{-(x+1)} + (e-1)e^{-ex} \quad x \in (0, \infty)$$

In this case, choosing $g(x) = e^{-x}$ gives

$$M = \sup_{x \in \mathbb{R}^+} \frac{f(x)}{g(x)} = \sup_{x \in \mathbb{R}^+} (e^{-1} + (e-1)e^{(1-e)x}) = e^{-1} + e - 1$$

If $f(x)$ is difficult to compute, then the test in step 2 can be slow to evaluate.

We can use simple **squeeze functions** that bracket f .

- For example, a squeeze function $s(y)$ below f has $s(y) \leq f(y)$.
- In a given trial, we evaluate $s(y)$ before $f(y)$. If $u \leq \frac{s(y)}{Mg(y)} \Rightarrow u \leq \frac{f(y)}{Mg(y)}$ and so we can automatically accept the point without computing $f(y)$.
- Clearly, we want the area under s to be as large as possible subject to $s(y) \leq f(y)$, so that we have a larger chance of accepting points directly.

Ratio of Uniforms

Theorem: Consider a distribution

$$f_X(x) = \frac{h(x)}{\int_{-\infty}^{\infty} h(x) dx} \quad \text{and} \quad \int_{-\infty}^{\infty} h(x) dx < \infty$$

Then for random variables (U, V) sampled uniformly distributed over the region

$$C_h = \{(u, v) : 0 \leq u \leq \sqrt{h(v/u)}\}$$

The random variable $X = V / U$ has PDF $f_X(x)$.

Proof: Let $Z = (U, V)$ be uniformly distributed on C_h .

$$f_Z(u, v) = \frac{1}{A} \mathbb{I}_{\{0 \leq u \leq \sqrt{h(v/u)}\}}$$

where A is the area of C_h . Now, let's apply the variable transformation:

$$\begin{aligned} X &= V / U & Y &= U \\ \Rightarrow U &= Y & V &= XY \end{aligned}$$

The transformation from $Z = (U, V)$ to $W = (X, Y)$

therefore has Jacobian

$$J = \begin{vmatrix} \partial u / \partial x & \partial u / \partial y \\ \partial v / \partial x & \partial v / \partial y \end{vmatrix} = \begin{vmatrix} 0 & 1 \\ y & x \end{vmatrix} = -y$$

The joint density of $W = (X, Y)$ is then given by

$$\begin{aligned} f_W(x, y) &= f_Z(u(x, y), v(x, y)) |J| \\ &= f_Z(y, xy) y \\ &= \frac{y}{A} \mathbb{I}_{\{0 \leq y \leq \sqrt{h(x)}\}} \end{aligned}$$

The marginal density of X is given by

$$f_X(x) = \int_0^{\sqrt{h(x)}} \frac{y}{A} dy = \left[\frac{y^2}{2A} \right]_0^{\sqrt{h(x)}} = \frac{h(x)}{2A}$$

However, f_X is a density, so

$$\begin{aligned} 1 &= \int_{-\infty}^{\infty} f_X(x) dx = \int_{-\infty}^{\infty} \frac{h(x)}{2A} dx \\ &\Rightarrow \int_{-\infty}^{\infty} h(x) dx = 2A \end{aligned}$$

So the marginal density of $X = V/U$ is given by

$$f_X(x) = \frac{h(x)}{\int_{-\infty}^{\infty} h(x) dx}$$

as desired. ■

This method is most useful if C_h can be contained in a rectangle. We therefore develop the following theorem:

Theorem: If $h(x)$ and $x^2 h(x)$ are bounded, then

$$C_h \subseteq [0, a] \times [b_-, b_+]$$

where

$$\begin{aligned} a &= \sup_{x \in \mathbb{R}} \sqrt{h(x)} \\ b_- &= -\sup_{x \leq 0} \sqrt{x^2 h(x)} \\ b_+ &= \sup_{x \geq 0} \sqrt{x^2 h(x)} \end{aligned}$$

Proof: If $(u, v) \in C_h$, then

$$0 < u < \sqrt{h(v/u)} \leq \sqrt{\sup_{x \in \mathbb{R}} h(x)} = \sup_{x \in \mathbb{R}} \sqrt{h(x)} = a$$

so u is indeed bounded as predicted.

For the v coordinate, we check two cases:

- If $\boxed{v \geq 0}$, then we substitute $t = \frac{v}{u} \geq 0$ into $0 \leq u^2 \leq h(\frac{v}{u})$ and get

$$0 \leq v^2 \leq t^2 h(t) \leq \sup_{t \geq 0} t^2 h(t) = b_+^2$$

Since v is positive, this implies $v < b_+$.

- If $\boxed{v < 0}$, a similar trick gives $v^2 < b_-^2$, and implies that $v > b_-$

We have therefore derived both bounds. ■

In the case where these conditions are met, we have a new method:

Ratio of uniforms: To generate a sample x from a distribution with CDF F

1. Simulate $u_1, u_2 \sim U(0,1)$
2. Let

$$u = au_1 \quad v = b_- + (b_+ - b_-)u_2$$
 This ensures that (u, v) is uniformly chosen from $[0, a] \times [b_-, b_+]$.
3. If $u^2 \leq h(\frac{v}{u}) \Rightarrow (u, v) \in C_h$, then return $x = v/u$.
Otherwise, return to step 1.

Analogously to rejection sampling, it may be possible to find sets $C_- \subseteq C_h \subseteq C_+$ for which it is easier to determine the membership of (u, v) . $[C_+ = [0, a] \times [b_-, b_+]$ is an example of such an “upper bounding” set].

EXAMPLE: Consider sampling from the Cauchy distribution

$$f(x) = \frac{1}{\pi(1 + x^2)}$$

Set

$$h(x) = \frac{1}{(1 + x^2)} \quad \text{clearly, } \int_{-\infty}^{\infty} h(x) \, dx = \pi$$

In this case, the exact form of C_h is

$$C_h = \{(u, v) : 0 \leq u, u^2 + v^2 \leq 1\}$$

$h(x)$ and $x^2 h(x)$ are both bounded [to see why, consider what happens as $x \rightarrow \infty$]. Now, we calculate

$$\begin{aligned}
 a &= \sup_{x \in \mathbb{R}^+} \sqrt{h(x)} = 1 \\
 b_- &= -\sup_{x \leq 0} \sqrt{x^2 h(x)} = -1 \\
 b_+ &= -\sup_{x \geq 0} \sqrt{x^2 h(x)} = 1
 \end{aligned}$$

As such, our “approximate” region is

$$C_h = \{(u, v) : u \in [0, 1], v \in [-1, 1]\}$$

Note that this method can also be used with $f(x) = h(x)$.

EXAMPLE: Consider sampling from the distribution

$$f(x) = e^{-(x+1)} + (e-1)e^{-ex} \quad x \in (0, \infty)$$

We use $h(x) = f(x)$. Now, the function is monotonously decreasing, and so reaches its maximum at $x = 0$ [**Note:** it is *important* to remember that this can be the case – differentiating and setting to 0 in such a case won’t help at all]. So

$$a = \sup_{x \in \mathbb{R}^+} \sqrt{f(x)} = \sqrt{e^{-1} + e - 1}$$

Similarly, the only non-positive value the function can have is 0, and so

$$b_- = -\sup_{x \leq 0} \sqrt{x^2 f(x)} = 0$$

Furthermore, by the triangle inequality $\sqrt{a+b} \leq \sqrt{a} + \sqrt{b}$, and so

$$b_+ = \sup_{x \geq 0} \sqrt{x^2 f(x)} \leq \sup_{x \geq 0} x \left\{ e^{-\frac{x-1}{2}} + (e-1)^{\frac{1}{2}} e^{-\frac{ex}{2}} \right\} = 2e^{-\frac{3}{2}} + \frac{2}{e^2} (e-1)^{\frac{1}{2}}$$

Composition

Let $\{f(x; \theta) : \theta \in \Theta\}$ denote a family of densities, and let $p(\theta)$ denote a density (or mass) over Θ . Then the density

$$f(x) = \begin{cases} \int_{\Theta} f(x; \theta) p(\theta) \, d\theta & \theta \text{ continuous} \\ \sum_{\theta \in \Theta} f(x; \theta) p(\theta) & \theta \text{ discrete} \end{cases}$$

is a **mixture density**. We can sample x from $f(x)$ as follows

- Generating θ with density/mass function $p(\theta)$
- Generating x with density $f(x; \theta)$

This is called the **method of composition**.

EXAMPLE: Consider sampling from the non-central chi-squared distribution, with $n \in \mathbb{N}$ degrees of freedom and non-centrality parameter $\lambda \in (0, \infty)$, denoted $\chi_n^2(\lambda)$

$$f(x; n, \lambda) = \sum_{r=0}^{\infty} \frac{e^{-\lambda/2} \lambda^r}{2^r r!} \frac{x^{\frac{n}{2}+r-1} e^{-\frac{x}{2}}}{2^{\frac{n}{2}+r} \Gamma\left(\frac{n}{2} + r\right)} \quad x \in (0, \infty)$$

Notice also that the distribution of the Poisson and chi-squared distributions are

$$f_{\text{Po}(\lambda)}(n; \lambda) = \frac{\lambda^n e^{-\lambda}}{n!} \quad f_{\chi_k^2}(x; k) = \frac{x^{\frac{k}{2}-1} e^{-\frac{x}{2}}}{2^{\frac{k}{2}} \Gamma\left(\frac{k}{2}\right)}$$

Clearly, therefore, we can generate the above by first generating $R \sim \text{Po}\left(\frac{\lambda}{2}\right)$, and then generating $X \sim \chi_{n+2R}^2$.

EXAMPLE: Consider sampling from the distribution

$$f(x) = e^{-(x+1)} + (e-1)e^{-ex} \quad x \in (0, \infty)$$

We note that this can be re-written as

$$f(x) = \frac{1}{e} e^{-x} + \frac{e-1}{e} \cdot e \cdot e^{-ex}$$

This is a mixture of exponentials, with

$$f(x; \theta) = \theta e^{-\theta x} \quad p(\theta) = \begin{cases} 1/e & \theta = 1 \\ (e-1)/e & \theta = e \end{cases}$$

We therefore generate two random variables, U_1 and U_2 , and

- If $U_1 > (1/e)$, return $-\log(U_2)$
- Else, return $-e^{-1} \log(U_2)$

Specific distributions

We now go through a number of commonly used algorithms for standard discrete and continuous distributions

- **Normal distribution** $Z \sim N(0,1)$
 - The **Box-Muller Method** uses $U_1, U_2 \sim U(0,1)$ and generates

$$Z_1 = \sqrt{-2 \log(U_1)} \cos(2\pi U_2) \quad Z_2 = \sqrt{-2 \log(U_1)} \sin(2\pi U_2)$$

This can be proven as follows

- The joint distribution of U_1 and U_2 is

$$f_{U_1, U_2}(u_1, u_2) = \mathbb{I}_{\{(u_1, u_2) \in (0,1)^2\}}$$

- $U_1 = \exp\left\{-\frac{1}{2}(Z_1^2 + Z_2^2)\right\}$ and $U_2 = \frac{1}{2\pi} \tan^{-1}(Z_2 / Z_1)$.
- We can then find the Jacobian

$$J = \begin{vmatrix} \partial U_1 / \partial Z_1 & \partial U_1 / \partial Z_2 \\ \partial U_2 / \partial Z_1 & \partial U_2 / \partial Z_2 \end{vmatrix}$$

- This allows us to write

$$\begin{aligned} f_{Z_1, Z_2}(z_1, z_2) &= f_{U_1, U_2}(u_1(z_1, z_2), u_2(z_1, z_2)) |J| \\ &= \frac{1}{2\pi} \exp\left\{-(z_1^2 + z_2^2) / 2\right\} \end{aligned}$$

Which is indeed the distribution of $N(0, 1)$.

- Evaluating trigonometric functions is slow. A faster method uses rejection sampling. Generate $V_1, V_2 \sim U(-1, 1)$ and set $r^2 = V_1^2 + V_2^2$. If $r^2 \geq 1$, reject. Otherwise, deliver

$$Z_1 = V_1 \sqrt{\frac{-2 \log(r^2)}{r^2}} \quad Z_2 = V_2 \sqrt{\frac{-2 \log(r^2)}{r^2}}$$

This can be proven as follows

- The random variables V_1 and V_2 that pass the rejection step are distributed uniformly on the unit disk, so $f_{V_1, V_2}(v_1, v_2) = \frac{1}{\pi} \mathbb{I}_{\{v_1^2 + v_2^2 \leq 1\}}$.

- We then write

$$V_1 = R \cos \Theta = \sqrt{R^2} \cos \Theta \quad V_2 = R \sin \Theta = \sqrt{R^2} \sin \Theta$$

- We find the Jacobian for the transformation involving R^2

$$J = \begin{vmatrix} \partial V_1 / \partial R^2 & \partial V_1 / \partial \Theta \\ \partial V_2 / \partial R^2 & \partial V_2 / \partial \Theta \end{vmatrix} = \begin{vmatrix} \frac{1}{2\sqrt{R^2}} \cos \Theta & -\sqrt{R^2} \sin \Theta \\ \frac{1}{2\sqrt{R^2}} \sin \Theta & \sqrt{R^2} \cos \Theta \end{vmatrix} = \frac{1}{2}$$

- This immediately gives

$$f_{R^2, \Theta}(r^2, \theta) = \frac{1}{2\pi} \mathbb{I}_{\{0 < R^2 \leq 1, 0 \leq \theta < \pi\}}$$

So R^2 and $\tilde{\Theta} = \frac{\Theta}{2\pi}$ are uniformly distributed on $(0, 1)$.

- We can therefore use the Box-Muller method with $U_1 = R^2$ and $U_2 = \tilde{\Theta}$. Since $\cos(2\pi\tilde{\Theta}) = \cos(\Theta) = V_1 / R$ and $\sin(2\pi\tilde{\Theta}) = V_2 / R$, the Box-Muller equations become exactly as above, in terms of R^2 , V_1 and V_2 .

- **Exponential distribution** $X \sim \exp(\mu)$: the inverse CDF method is easy to implement and considered satisfactory. Note that the inverse CDF is $F^{-1}(U) = -\mu^{-1} \log(U)$.

- **Gamma distribution** $X \sim \Gamma(\alpha, \beta)$: Note that if $Z \sim \Gamma(\alpha, 1)$, then $X = Z / \beta$, so sampling $Z \sim \Gamma(\alpha, 1)$ is enough. A number of algorithms based on rejection methods exist depending on whether $\alpha > 1$ or $\alpha < 1$. (If $\alpha = 1$, we have an exponential distribution).

For an inverse gamma, note that $X \sim \Gamma^{-1}(\alpha, \beta) \Rightarrow X^{-1} \sim \Gamma(\alpha, \beta)$, so simply take the reciprocal of gamma draws.

- **Chi-squared distribution** $X \sim \chi_\nu^2$: This is a special case of the gamma; $X \sim \Gamma(\frac{\nu}{2}, \frac{1}{2})$.

For small ν , we can use the Box-Muller method to generate ν normally distributed variables, square them, and add them together. Note that if Z_1 and Z_2 are two Box-Muller generated variables, $Z_1^2 + Z_2^2 = -2\log(U_1)$, and so when ν is even, the sum is simply given by $-2\log(\prod_{i=1}^{\nu/2} U_i)$.

- **Poisson distribution** $X \sim \text{Po}(\lambda)$: A slow but clever way of sampling X is to realise that if the number of arrivals in an interval $[0, t]$ is Poisson distributed with mean λt , then the time between each Poisson arrival is distributed as $\exp(\lambda^{-1})$, for which $F^{-1}(U) = -\frac{1}{\lambda}\log(U)$.

Thus, we simply generate such exponential variables and continuously sum them until the sum is greater than 1 – let N be the number of random variables required for this to happen. Then N is a realisation of X .

Non-parameteric Inference

Given a sample X_1, \dots, X_n of independent variables with distribution F , we are often interested in estimating some parameter $\theta = \theta(F)$. Common examples:

- $\theta = \mathbb{E}_f \{ \phi(x) \}$
- $\theta = P_f (X_1 \in A)$ – this is a special case of the above with $\phi(x) = \mathbb{I}_{\{x \in A\}}$.
- $\theta = F^{-1}(1/2)$, the median.

The **plugin principle** is often used in estimating such values:

Plugin estimator: The plugin estimator for θ is

$$\hat{\theta} = \theta(\hat{F})$$

Where \hat{F} is an estimator of F . Often, the **empirical distribution function** (ECDF) \hat{F}_n is used:

$$\hat{F}_n(x) = \frac{1}{n} \sum_{i=1}^n \mathbb{I}_{\{x_i \leq x\}}$$

in which case we denote the plugin estimate $\hat{\theta}_n$.

R-CODE: The ECDF of a set of points in R can be found using

```
ecdf(dataPoints)
```

So for example, to plot find the ECDF of the normal distribution, based on 10 points, use

```
plot(ecdf(rnorm(10, 0, 1)), verticals=TRUE)
```

(The last statement ensures vertical lines are drawn to connect the “steps” in the graph).

To then add a line representing the *real* density, use

```
x <- seq(-2.5, 2.5, length=1000)
```

(this generates 1000 x values between -2.5 and 2.5), and then

```
lines(x, pnorm(x, 0, 1), col=2, lty=2)
```

(the last two statements make the line red and dotted).

Using the ECDF, the plugin estimators of the above examples are

- $\hat{\theta}_n = \frac{1}{n} \sum \phi(x_i)$
- $\hat{\theta}_n = \frac{1}{n} \sum \mathbb{I}_{\{x_i \in A\}}$
- $\hat{\theta}_n = F_n^{-1}(\frac{1}{2}) = x_{(\lfloor n/2 \rfloor)}$, where $X_{(1)} < \dots < X_{(n)}$ are the ordered X .

These estimators are unbiased, and it can be shown³ that

$$\text{Var}(\hat{\theta}_n) = \frac{1}{n} \text{Var}\{\phi(X)\}$$

This variance, however, may be very large. In **variance reduction**, we try to reduce the variance of estimators while maintaining other good qualities.

Importance Sampling (IS)

Importance sampling reduces the variance of our estimator, and also allows us to sample from another, simpler distribution. We define the **support** of a function g , \mathcal{Y}_g as $\mathcal{Y}_g = \{y : g(y) > 0\}$.

Importance sampling: Suppose we are trying to estimate $\theta = \mathbb{E}_f\{\phi(x)\}$. The density of interest is f , and let g denote another density that is easily sampled from and such that

$$f(x)|\phi(x)| > 0 \Rightarrow g(x) > 0$$

(Or in other words, $\mathcal{Y}_{f(x)|\phi(x)} \subseteq \mathcal{Y}_g$).

Sample y_1, \dots, y_n independently from g and consider estimators of the form

$$\hat{\theta}_g = \frac{1}{n} \sum_{i=1}^n w_i \phi(y_i) \quad w_i = w(y_i) = \frac{f(y_i)}{g(y_i)}$$

where the w_i are called **importance weights**.

Then the estimator $\hat{\theta}_g$ is also unbiased, and its variance is minimized when

$$g(x) = g_0(x) = \frac{|\phi(x)| f(x)}{\int_{\mathcal{Y}_{f|\phi}} |\phi(y)| f(y) \, dy}$$

³ We prove this result as follows

$$\text{Var}\left(\hat{\theta}_n\right) = \text{Var}\left(\frac{1}{n} \sum \phi(x_i)\right) = \frac{1}{n^2} \text{Var}\left(\sum \phi(x_i)\right) = \frac{1}{n} \text{Var}\{\phi(X)\}$$

(note that in the last step, we only extra a factor of n from the variance – not n^2 – because the different items in the sum are independent).

Intuitively – the last line shows that we are sampling from a region where both f and ϕ are large; in other words, very informative regions.

Proof: We first show that our estimator is unbiased

$$\begin{aligned}
 \mathbb{E}_g \{ \hat{\theta}_g \} &= \frac{1}{n} \sum_{i=1}^n \mathbb{E} \{ w(Y) \phi(Y) \} \\
 &= \mathbb{E}_g \{ w(Y) \phi(Y) \} \\
 &= \int_{\mathcal{Y}_g} \left(\frac{f(y)}{g(y)} \phi(y) \right) g(y) \, dy \\
 &= \int_{\mathcal{Y}_g} \phi(y) f(y) \, dy \\
 &= \int_{\mathcal{Y}_{f|\phi}} \phi(y) f(y) \, dy \\
 &= \mathbb{E}_f \{ \phi(x) \} \\
 &= \theta
 \end{aligned}$$

(We were able to change the range of integration in the second-to-last step, because in all those members of \mathcal{Y}_g that are not members of $\mathcal{Y}_{f|\phi}$, the integral would have been 0 anyway).

The variance of our estimator is given by

$$\begin{aligned}
 \text{var}_g \{ \hat{\theta}_g \} &= \text{var}_g \left\{ \frac{1}{n} \sum_{i=1}^n w(Y) \phi(Y) \right\} \\
 &= \frac{1}{n} \text{var}_g \{ w(Y) \phi(Y) \} \\
 &= \frac{1}{n} \text{var}_g \left\{ \frac{f(Y)}{g(Y)} \phi(Y) \right\}
 \end{aligned}$$

So

$$\begin{aligned}
 n \text{var}_g \{ \hat{\theta}_g \} &= \mathbb{E}_g \left\{ \left(\frac{f(Y)}{g(Y)} \phi(Y) \right)^2 \right\} - \mathbb{E}_g \left\{ \frac{f(Y)}{g(Y)} \phi(Y) \right\}^2 \\
 &= \int_{\mathcal{Y}_g} \frac{f(y)^2 \phi(y)^2}{g(y)^2} g(y) \, dy - \theta^2 \\
 &= \int_{\mathcal{Y}_{f|\phi}} \frac{f(y)^2 \phi(y)^2}{g(y)} \, dy - \theta^2
 \end{aligned}$$

(We were able to change the range of integration in the last step because at any point in \mathcal{Y}_g which is not in $\mathcal{Y}_{f|\phi}$, the integral is 0 anyway).

Clearly, the variance is minimized when the integral in the last line above is minimized. By Jensen's Inequality, however, we can find this lower bound:

$$\begin{aligned}\mathbb{E}_g \left\{ \left(\frac{f(Y)}{g(Y)} \phi(Y) \right)^2 \right\} &\geq \mathbb{E}_g \left\{ \left| \frac{f(Y)\phi(Y)}{g(Y)} \right|^2 \right\} \\ &= \left(\int_{\mathcal{Y}_{f|\phi}} |\phi(Y)| f(y) \, dy \right)^2\end{aligned}$$

This bound is obviously achieved when g is as given in the theorem. ■

The only issue is that if we know the integral in the expression for g_0 , we probably know the integral of interest anyway! It is furthermore unclear how we might sample from g_0 . The theorem is useful, however, in suggesting that we should seek a g “close” to g_0 from which it is easy to sample.

Another interesting note: if $\phi(x) = \mathbb{I}_{\{x \in A\}}$, then g_0 is simply the conditional density of X given $x \in A$... This makes sense; this is where the data is most informative (we don't care about points outside A).

A final point is that if f or g are only known up to a normalising constant, then the method can still be used with the estimator

$$\hat{\theta}_g = \frac{\sum_{i=1}^n w_i \phi(x_i)}{\sum_{i=1}^n w_i}$$

is asymptotically unbiased, because

$$\mathbb{E}_g(W) = \mathbb{E}_g \{ f(Y) / g(Y) \} = 1$$

so as $n \rightarrow \infty$, $\frac{1}{n} \sum w_i \rightarrow 1$. Short of such an asymptotic situation, the estimator above will always exhibit some (small) bias. However, choosing g correctly will still give us a small variance estimator.

This is the **bias-variance tradeoff**.

This method can be used to estimate the CDF F by allowing ϕ to take two arguments.

$$\hat{\theta}_g(x) = \frac{\sum_{i=1}^n w_i \phi(y_i, x)}{\sum_{i=1}^n w_i} = \sum_{i=1}^n \frac{w_i}{\sum_{i=1}^n w_i} \mathbb{I}_{\{y_i \leq x\}} = \hat{F}_{g,n}(x)$$

This can easily be worked out as follows

- Sample y_i from g [a uniform distribution, for example]. Let $y_{(i)}$ be the i^{th} ordered part of the sample.
- Work out the weights, normalise them – call the normalised weights \tilde{w} .
- Let $\hat{F}_{g,n}(y_{(1)}) = \tilde{w}(y_{(1)})$
- Thence, $\hat{F}_{g,n}(y_{(i+1)}) = \hat{F}_{g,n}(y_{(i)}) + \tilde{w}(y_{(i+1)})$

Sometimes, it is useful to retain a sample from f rather than an estimate of a parameter. This can be achieved by **sampling with replacement** from y_1, \dots, y_n with a discrete distribution proportional to w_1, \dots, w_n . This is called **sampling importance re-sampling** (SIR).

R-CODE: The following function

```
sample(y, num, replace=TRUE, prob=w)
```

Samples num items from the vector y , *with* replacement, and places a probability w on each item.

EXAMPLE: Suppose that X has Cauchy distribution, with density

$$f(x) = \frac{1}{\pi(1+x^2)} \quad x \in \mathbb{R}$$

and cumulative density

$$F(x) = \frac{1}{\pi} \arctan(x) + \frac{1}{2}$$

and we want to estimate $\theta = \mathbb{P}(X > 2)$. In other words, we have

$$\phi(x) = \mathbb{I}_{\{x > 2\}}.$$

We want g to be as close as possible to

$$\begin{aligned}
 f(X = x | X > 2) &= \frac{f(x > 2 | X = x)f(X = x)}{\mathbb{P}(X > 2)} \\
 &= \frac{\mathbb{I}_{\{x>2\}} \frac{1}{\pi(1+x^2)}}{1 - \frac{1}{\pi} \arctan(2) + \frac{1}{2}} \\
 &= \frac{1}{\left(\frac{\pi}{2} - \arctan 2\right)(1+x^2)} \mathbb{I}_{\{x>2\}}
 \end{aligned}$$

A sensible choice seems to be $g(x) = \frac{2}{x^2} \mathbb{I}_{\{x>2\}}$, which is easy to sample from by inversion [the factor of 2 is to ensure the distribution is normalised].

We then sample y_i from g . In this case, every value we sample is greater than 2, and so $\phi(x) = 1$. As such, the IS estimator is simply

$$\hat{\theta}_g = \frac{1}{n} \sum_{i=1}^n w_i \quad w_i = \frac{f(y_i)}{g(y_i)} = \frac{y_i^2}{2\pi(1+y_i^2)}$$

Consider the variance of this estimator

$$\begin{aligned}
 \text{var}_g \{ \hat{\theta}_g \} &= \text{var}_g \left\{ \frac{1}{n} \sum_{i=1}^n \frac{y_i^2}{2\pi(1+y_i^2)} \right\} \\
 &= \frac{1}{4\pi^2 n} \text{var}_g \left\{ \frac{y_i^2}{1+y_i^2} \right\} \\
 &= \frac{1}{4\pi^2 n} \left\{ \int_{-\infty}^{\infty} \frac{y_i^4}{(1+y_i^2)^2} \frac{2}{y_i^2} \mathbb{I}_{\{x>2\}} dx - \theta^2 \right\} \\
 &= \frac{1}{2\pi^2 n} \left\{ \int_2^{\infty} \frac{y_i^2}{(1+y_i^2)^2} dx - \theta^2 \right\} \\
 &= \frac{1}{n} \left(\frac{1}{10\pi^2} + \frac{\theta}{4\pi} - \theta^2 \right)
 \end{aligned}$$

Furthermore,

$$\begin{aligned}
 \text{var}_f \{ \hat{\theta}_n \} &= \text{var}_f \left\{ \frac{1}{n} \sum_{i=1}^n \mathbb{I}_{\{x_i>2\}} \right\} \\
 &= \frac{1}{n} \text{var}_f \{ \mathbb{I}_{\{X>2\}} \} \\
 &= \frac{1}{n} \left\{ [0 \cdot \mathbb{P}(X < 2) + 1 \cdot \mathbb{P}(X > 2)] - \theta^2 \right\} \\
 &= \frac{1}{n} \theta(1 - \theta)
 \end{aligned}$$

Feeding in the true value of $\theta = \frac{1}{2} - \frac{1}{\pi} \arctan 2$, we find that the independence sampler is a large improvement.

Control variates

Definition (Control variates): Suppose Y is an unbiased estimator of θ . C is a *control variate* for Y if it is correlated with Y and its mean μ_C is known.

Control variates: Suppose Y is an unbiased estimator of θ , and C is a *control variate* for Y . Then

$$Y_\beta = Y - \beta(C - \mu_C)$$

is also unbiased for θ , and its variance is minimized when

$$\beta = \beta^* = \frac{\text{cov}\{Y, C\}}{\text{var}\{C\}}$$

at which point

$$\text{var}\{Y_{\beta^*}\} = (1 - \rho^2) \text{var}\{Y\} \leq \text{var}\{Y\}$$

where $\rho = \text{corr}\{Y, C\}$.

Proof: It is pretty obvious that Y_β is unbiased. We then have

$$\begin{aligned} \text{var}\{Y_\beta\} &= \text{var}\{Y - \beta C + \beta\mu_C\} \\ &= \text{var}\{Y - \beta C\} \\ &= \mathbb{E}\left\{\left(\{Y - \theta\} + \{-\beta C + \beta\mu_C\}\right)^2\right\} \\ &= \mathbb{E}\left\{(Y - \theta)^2\right\} + \mathbb{E}\left\{\beta^2 (C - \mu_C)^2\right\} \\ &\quad - 2\beta \mathbb{E}\left\{\{Y - \theta\}\{C - \mu_C\}\right\} \\ &= \text{var}\{Y\} + \beta^2 \text{var}\{C\} - 2\beta \text{cov}\{Y, C\} \end{aligned}$$

Differentiating with respect to β and setting to 0, we obtain the results above. ■

Note that even if $\text{cov}(Y, C)$ is not known, there is an optimal β that relies only on the size of the covariance. Furthermore, the ideas can be extended to more than one control variate as follows:

$$\tilde{Y} = Y - \beta_1(C_1 - \mu_{C_1}) - \cdots - \beta_k(C_k - \mu_{C_k})$$

EXAMPLE: We return to the Cauchy example of the last section, in which we were estimating $\theta = \mathbb{P}(X > 2)$. A different approach would be to use the estimator

$$\theta = \frac{1}{2} - \mathbb{P}(0 < X < 2) = \frac{1}{2} - \int_0^2 f(x) \, dx$$

We can estimate the integral using Monte-Carlo integration. In this case, we use $\phi(x) = \mathbb{I}_{\{x \in (0,2)\}}$, and

$$g \sim U(0,2) \Rightarrow g(x) = \frac{1}{2} \mathbb{I}_{\{x \in (0,2)\}}$$

We then generate y_1, \dots, y_n independently of g . Every one of these values will be between 0 and 2, and so $\phi(x) = 1$. The IS estimator of the integral based on g is then

$$\frac{1}{n} \sum_{i=1}^n \frac{f(y_i)}{\frac{1}{2}} = \frac{2}{n} \sum_{i=1}^n f(y_i)$$

and so our estimator of θ will be

$$\tilde{\theta} = \frac{1}{2} - \frac{2}{n} \sum_{i=1}^n f(y_i)$$

A Taylor Expansion of f suggests the following improved estimator:

$$\tilde{\theta} = \frac{1}{2} - \frac{2}{n} \sum_{i=1}^n f(y_i) - \beta (y_i^2 - \mathbb{E}[y_i^2])$$

This is clearly of control variate form, with

$$C = \frac{2}{n} \sum_{i=1}^n y_i^2$$

The variance of this estimator is minimised at

$$\beta = \beta^* = \frac{\text{cov}\{Y, C\}}{\text{var}\{C\}}$$

Now remember that the different y_i are independent. Thus

$$\begin{aligned} \text{cov}\{Y, C\} &= \frac{4}{n^2} \text{cov}\left\{\sum f(y_i), \sum y_i^2\right\} \\ &= \frac{4}{n^2} n \text{cov}\left\{\frac{1}{\pi(1+Y^2)}, Y^2\right\} \\ &= \frac{4}{n} \left\{ \mathbb{E}\left[\frac{Y^2}{\pi(1+Y^2)}\right] - \mathbb{E}\left[\frac{1}{\pi(1+Y^2)}\right] \mathbb{E}(Y^2) \right\} \end{aligned}$$

Remembering Y is uniformly distributed, the above are relatively simple to calculate. Similarly

$$\text{var} \left\{ \frac{2}{n} \sum_{i=1}^n y_i^2 \right\} = \frac{4}{n} \text{var} \{ Y^2 \}$$

We find

$$\frac{\text{cov} \{ Y, C \}}{\text{var} \{ C \}} = \frac{\text{cov} \left\{ \frac{1}{\pi(1+Y^2)}, Y^2 \right\}}{\text{var} \{ Y^2 \}} = \frac{45}{64\pi} \left(1 - \frac{7}{6} \arctan 2 \right)$$

Finding the actual variance of the estimator involves similar steps.

Antithetic variables

Antithetic variables: Suppose $\hat{\theta}_1$ and $\hat{\theta}_2$ are two estimators of θ with the same expectation. We say they are *antithetic variables* if they are negatively correlated.

We then consider estimators formed as a convex combination of $\hat{\theta}_1$ and $\hat{\theta}_2$

$$\hat{\theta}_\lambda = \lambda \hat{\theta}_1 + (1 - \lambda) \hat{\theta}_2 \quad \lambda \in (0, 1)$$

We have $\mathbb{E}(\hat{\theta}_\lambda) = \mathbb{E}(\hat{\theta}_1) = \mathbb{E}(\hat{\theta}_2)$. The variance of the estimator is minimized by

$$\lambda^* = \frac{\sigma_2^2 - \rho \sigma_1 \sigma_2}{\sigma_1^2 + \sigma_2^2 - 2\rho \sigma_1 \sigma_2}$$

at which point it takes value

$$\text{var}(\hat{\theta}_{\lambda^*}) = \frac{\sigma_1^2 \sigma_2^2 (1 - \rho^2)}{\sigma_1^2 + \sigma_2^2 - 2\rho \sigma_1 \sigma_2}$$

Where $\sigma_1^2 = \text{var}(\hat{\theta}_1)$, $\sigma_2^2 = \text{var}(\hat{\theta}_2)$ and $\rho = \text{corr}(\hat{\theta}_1, \hat{\theta}_2)$.

Proof: We have

$$\begin{aligned} \text{var}(\hat{\theta}_\lambda) &= \text{var} \left\{ \lambda \hat{\theta}_1 + (1 - \lambda) \hat{\theta}_2 \right\} \\ &= \mathbb{E} \left\{ \left[\lambda (\hat{\theta}_1 - \mu_1) + (1 - \lambda) (\hat{\theta}_2 - \mu_2) \right]^2 \right\} \\ &= \lambda^2 \mathbb{E} \left\{ (\hat{\theta}_1 - \mu_1)^2 \right\} + (1 - \lambda)^2 \mathbb{E} \left\{ (\hat{\theta}_2 - \mu_2)^2 \right\} \\ &\quad + 2\lambda(1 - \lambda) \mathbb{E} \left\{ (\hat{\theta}_1 - \mu_1) (\hat{\theta}_2 - \mu_2) \right\} \\ &= \lambda^2 \sigma_1^2 + (1 - \lambda)^2 \sigma_2^2 + 2\lambda(1 - \lambda) \text{cov}(\hat{\theta}_1, \hat{\theta}_2) \\ &= \lambda^2 \sigma_1^2 + (1 - \lambda)^2 \sigma_2^2 + 2\lambda(1 - \lambda) \rho \sigma_1 \sigma_2 \end{aligned}$$

Differentiating with respect to λ and setting to 0, we obtain the results above. ■

Note: the computational cost of computing $\hat{\theta}_{\lambda^*}$ is twice that of computing each of the individual estimators, but we get a significant reduction in variance.

The method is especially useful when used with the following theorem. (Note: a **non-degenerate distribution** is one for which there is no $a \in \mathbb{R}$ such that $F(x) = \mathbb{I}_{\{x \geq a\}}$)

Theorem: Let F be a non-degenerate distribution function and $U \sim U(0,1]$, then

$$\text{cov}\{F^{-1}(U), F^{-1}(1-U)\} < 0$$

Proof: Let

$$\theta = \mathbb{E}\{F^{-1}(U)\} = \mathbb{E}\{F^{-1}(1-U)\}$$

Observe that

$$F^{-1}(1-u) > \theta \Leftrightarrow u < 1 - F(\theta)$$

We have

$$\begin{aligned} & \text{cov}\{F^{-1}(U), F^{-1}(1-U)\} \\ &= \mathbb{E}\left\{\left(F^{-1}(U) - \theta\right)\left(F^{-1}(1-U) - \theta\right)\right\} \\ &= \mathbb{E}\left\{F^{-1}(U)\left(F^{-1}(1-U) - \theta\right)\right\} \\ &= \int_0^{1-F(\theta)} F^{-1}(u)\left[F^{-1}(1-u) - \theta\right] du \\ &\quad + \int_{1-F(\theta)}^1 F^{-1}(u)\left[F^{-1}(1-u) - \theta\right] du \end{aligned}$$

We want to try and get an upper bound for this quantity; ie: find the **largest** it could ever be. To do that, we note that

- The quantity in the square brackets is positive in the first integral.
- Thus, if we replace $F^{-1}(u)$ by a constant equal to the largest value it can take in that range of integration, we get an upper bound.

- Since $F^{-1}(u)$ is a non-decreasing function, the largest value it can take in the range of the first integral is $F^{-1}(1 - F(\theta))$.
- So the first integral is smaller than

$$F^{-1}(1 - F(\theta)) \int_0^{1-F(\theta)} F^{-1}(1 - u) - \theta \, du$$

The quantity in the square brackets is negative in the second integral, and a similar argument applies.

Together, these imply that

$$\begin{aligned} & \text{cov}\{F^{-1}(U), F^{-1}(1 - U)\} \\ & < F^{-1}(1 - F(\theta)) \left\{ \int_0^{1-F(\theta)} F^{-1}(1 - u) - \theta \, du \right. \\ & \quad \left. + \int_{1-F(\theta)}^1 F^{-1}(1 - u) - \theta \, du \right\} \\ & = F^{-1}(1 - F(\theta)) \int_0^1 F^{-1}(1 - u) - \theta \, du \\ & = F^{-1}(1 - F(\theta)) \left\{ \int_0^1 F^{-1}(1 - u) \, du - \theta \right\} \\ & = 0 \end{aligned}$$

Our covariance is therefore negative. ■

Effectively, the method of inversion provides two correlated samples; one based on U and one based on $1 - U$. This method is useful when we have an estimator of the form $\hat{\theta} = \frac{1}{n} \sum_{i=1}^n g(y_i)$, and the y_i can be simulated by inversion.

EXAMPLE: Back once again to the Cauchy distribution, we return to our estimator of the form

$$\tilde{\theta} = \frac{1}{2} - \frac{2}{n} \sum_{i=1}^n \frac{1}{\pi(1 + y_i^2)}$$

where $y_1, \dots, y_n \sim U(0,2)$, IID. An antithetic variable estimator based on this is simply

$$\tilde{\theta} = \frac{1}{2} - \frac{2}{n} \sum_{i=1}^n \frac{1}{2 \pi(1 + y_i^2)} + \frac{1}{2} \frac{1}{\pi(1 + (2 - y_i)^2)}$$

EXAMPLE: Consider that

$$\int_0^1 \sqrt{1 - u^2} \, du = \frac{\pi}{4}$$

This means that a Monte-Carlo estimator of π is

$$\hat{\pi} = \frac{1}{n} \sum_{i=1}^n 4\sqrt{1-U_i^2}$$

Where U_i are independent $U(0,1)$ random variables. An estimator based on the method of antithetic variables is

$$\hat{\pi}_{AV} = \frac{1}{n} \sum_{i=1}^n 4 \left\{ \frac{1}{2} \sqrt{1-U_i^2} + \frac{1}{2} \sqrt{1-(1-U_i)^2} \right\}$$

Now, note that

$$\begin{aligned} \text{Var}(\hat{\pi}) &= \frac{1}{n^2} \text{Var} \left(\sum_{i=1}^n 4\sqrt{1-U_i^2} \right) \\ &= \frac{16}{n} \text{Var} \left(\sqrt{1-U_i^2} \right) \end{aligned}$$

And

$$\begin{aligned} \text{Var}(\hat{\pi}_{AV}) &= \frac{1}{n^2} \text{Var} \left[\sum_{i=1}^n 4 \left\{ \frac{1}{2} \sqrt{1-U_i^2} + \frac{1}{2} \sqrt{1-(1-U_i)^2} \right\} \right] \\ &= \frac{4}{n} \left[\text{Var} \left(\sqrt{1-U_i^2} \right) + \text{Var} \left(\sqrt{1-(1-U_i)^2} \right) \right. \\ &\quad \left. + 2\text{Cov} \left(\sqrt{1-U_i^2}, \sqrt{1-(1-U_i)^2} \right) \right] \\ &= \frac{4}{n} \left[2\text{Var} \left(\sqrt{1-U_i^2} \right) + 2\text{Cov} \left(\sqrt{1-U_i^2}, \sqrt{1-(1-U_i)^2} \right) \right] \end{aligned}$$

Note that there is no factor of n in front of the covariance because the U_i are independent and therefore uncorrelated with each other. Thus, cross terms vanish from the covariance.

As such

$$\begin{aligned} \frac{\text{Var}(\hat{\pi}_{AV})}{\text{Var}(\hat{\pi})} &= \frac{1}{4} \frac{\left[2\text{Var} \left(\sqrt{1-U_i^2} \right) + 2\text{Cov} \left(\sqrt{1-U_i^2}, \sqrt{1-(1-U_i)^2} \right) \right]}{\text{Var} \left(\sqrt{1-U_i^2} \right)} \\ &= \frac{1}{2} \left\{ 1 + \frac{\text{Cov} \left(\sqrt{1-U_i^2}, \sqrt{1-(1-U_i)^2} \right)}{\text{Var} \left(\sqrt{1-U_i^2} \right)} \right\} \\ &= \frac{1}{2} \left\{ 1 + \frac{\int_0^1 \sqrt{1-u^2} \sqrt{1-(1-u)^2} \, du - \frac{\pi^2}{16}}{\left(1 - \frac{1}{3}\right) - \frac{\pi^2}{16}} \right\} \\ &= 0.140 \end{aligned}$$

Non-parameteric Inference

Monte-Carlo Tests

Let X_1, \dots, X_n be independent with distribution function F , and suppose we want to use a statistic $T = T(X_1, \dots, X_n)$ to test

$$\begin{aligned} H_0 : F &= F_0 && \text{against} \\ H_1 : F &\neq F_0 \end{aligned}$$

If **small values of T** represent **departure from H_0** , then a test of size $\alpha \in (0,1)$ would **reject H_0** if $T < c_\alpha$, where c_α is the α^{th} quantile of T .

If the null distribution of T is unknown, however, we may not be able to compute c_α . In a **Monte Carlo test** of approximate size α , we estimate c_α as follows:

1. Choose a large $B \in \mathbb{N}$.
2. For each $k \in \{1, \dots, B\}$, associate an α , namely $\alpha = k / (B + 1)$, and restrict the choice of α to those values.
3. For each $k \in \{1, \dots, B\}$, assume H_0 is true, and simulate a random sample of N variables

$$X_{k1}^*, \dots, X_{kn}^* \quad \text{for } k = 1, \dots, B$$
4. For each $k \in \{1, \dots, B\}$, compute $T_k^* = T(X_{k1}^*, \dots, X_{kn}^*)$
5. Let $c_\alpha \equiv T_{(k)}^*$, where $T_{(k)}^*$ is the k^{th} item when the T_k^* are ordered.

Since the critical point is random, the critical region is “blurred”

Theorem: Assume that under H_0 , T has density f_0 supported on an interval. Then the MC test has *exact* size α .

Proof: The size of the test is (by definition) $\mathbb{P}(T < T_{(k)}^* \mid H_0 \text{ true})$. We can condition on the value of T :

$$\begin{aligned} \mathbb{P}(T < T_{(k)}^*) &= \int_{-\infty}^{\infty} \mathbb{P}(T < T_{(k)}^* \mid T = t) f_0(t) dt \\ &= \int_{-\infty}^{\infty} \mathbb{P}(T_{(k)}^* > t) f_0(t) dt \end{aligned}$$

Note, however, that $\mathbb{P}(T_{(k)}^* > t)$ is just the probability that *less than* k of the T^* are $\leq t$, so

$$\begin{aligned} \mathbb{P}(T_{(k)}^* > t) &= \mathbb{P}\left(\text{bin}(B, \mathbb{P}(T^* \leq t)) < k\right) \\ &= \sum_{r=0}^{k-1} {}^B C_r F_0(t)^r \{1 - F_0(t)\}^{B-r} \end{aligned}$$

Note that this last line assumes the T^* have the same distribution as T . It is therefore *not* acceptable to obtain the X^* using the jackknife or bootstrap.

Feeding this back into our integral

$$\begin{aligned} \mathbb{P}(T < T_{(k)}^*) &= \int_{-\infty}^{\infty} \sum_{r=0}^{k-1} {}^B C_r F_0(t)^r \{1 - F_0(t)\}^{B-r} f_0(t) dt \\ &= \int_{-\infty}^{\infty} \sum_{r=0}^{k-1} {}^B C_r u^r \{1 - u\}^{B-r} du \end{aligned}$$

(The last line uses the existence of a density). Since the sum is finite:

$$\begin{aligned} &= \sum_{r=0}^{k-1} \int_{-\infty}^{\infty} {}^B C_r u^r \{1 - u\}^{B-r} du \\ &= \sum_{r=0}^{k-1} r^{-1} \int_{-\infty}^{\infty} u r {}^B C_r u^{r-1} \{1 - u\}^{B-r} du \end{aligned}$$

Replacing the factorials in the binomial coefficients with gamma functions, we see that this equivalent to taking $\mathbb{E}\{\beta(r, B - r + 1)\} = r / (B + 1)$, which gives

$$\begin{aligned} &= \sum_{r=0}^{k-1} r^{-1} \frac{r}{B + 1} \\ &= \frac{k}{B + 1} \\ &= \alpha \end{aligned}$$

As required. ■

The hard part of the process above is step 3; simulating random samples under H_0 . This can be done using **resampling methods**, which involve the use of samples taken from a single observed sample, and can be used when very little is known about the underlying distribution. We study two such methods...

The Jackknife

The jackknife involves subsampling without replacement from an observed sample. It involves the use of leave-one-out data sets and is most often used to estimate such quantities as the variance or bias of an estimator.

We first develop some notation. Let X_1, \dots, X_n be a sample of independent variables with distribution F , and let $\theta = \theta(F)$ be a parameter of interest. Let $\hat{\theta}^n = \hat{\theta}^n(X_1, \dots, X_n)$ be an estimator of θ with variance v , and for $i = 1, \dots, n$ let $\hat{\theta}_{(-i)}^{n-1} = \hat{\theta}^{n-1}(X_1, \dots, X_{i-1}, X_{i+1}, \dots, X_n)$ be an estimator based on all the variables in the sample save one.

Jackknife estimator of variance: The *jackknife estimator of variance* is given by

$$\hat{v}_{\text{jack}} = \frac{n-1}{n} \sum_{i=1}^n \left(\hat{\theta}_{(-i)}^{n-1} - \hat{\theta}_{(\text{av})}^{n-1} \right)^2$$

where

$$\hat{\theta}_{(\text{av})}^{n-1} = \frac{1}{n} \sum_{i=1}^n \hat{\theta}_{(-i)}^{n-1}$$

Motivation: We can motivate this definition by considering a **linear statistic** of the form

$$\hat{\theta} = \mu + \frac{1}{n} \sum_{i=1}^n \alpha(X_i)$$

for some function α .

In this case

$$\hat{\theta}_{(-i)}^{n-1} = \mu + \frac{1}{n-1} \sum_{j \neq i} \alpha(X_j)$$

and so

$$\begin{aligned} \hat{\theta}_{(\text{av})}^{n-1} &= \mu + \frac{1}{n} \sum_{i=1}^n \frac{1}{n-1} \sum_{j \neq i} \alpha(X_j) \\ &= \mu + \frac{1}{n} \frac{1}{n-1} \sum_{i=1}^n \sum_{j \neq i} \alpha(X_j) \end{aligned}$$

note that the sum involves every $\alpha(X_i)$ exactly $n-1$ times, because we sum from $i=1 \rightarrow n$, leaving out one item each time. We can therefore write it as

$$\begin{aligned}
&= \mu + \frac{1}{n} \sum_{i=1}^n \alpha(X_i) \\
&\equiv \mu + \bar{\alpha}
\end{aligned}$$

Where $\bar{\alpha}$ denotes the average of $\alpha(X_i)$.

We can now write:

$$\begin{aligned}
\hat{v}_{\text{jack}} &= \frac{n-1}{n} \sum_{i=1}^n \left(\hat{\theta}_{(-1)}^{n-1} - \mu - \bar{\alpha} \right)^2 \\
&= \frac{n-1}{n} \sum_{i=1}^n \left(\frac{1}{n-1} \sum_{j \neq i} \alpha(X_j) - \bar{\alpha} \right)^2 \\
&= \frac{n-1}{n} \sum_{i=1}^n \left(\frac{1}{(n-1)^2} \sum_{j,k \neq i} \alpha(X_j) \alpha(X_k) \right. \\
&\quad \left. - \frac{2\bar{\alpha}}{n-1} \sum_{j \neq i} \alpha(X_j) + \bar{\alpha}^2 \right)^2
\end{aligned}$$

We now combine the outer and inner sums as follows:

- The first sum is hard! We can decompose it into two parts:
 - First, we'll get terms of the form $\alpha(X_i)^2$. We'll get $n-1$ of each of those, because we sum from $i=1 \rightarrow n$, leaving out one item each time.
 - We'll then get terms of the form $\alpha(X_i)\alpha(X_j)$. We'll get $n-2$ of each one of those, because we sum from $i=1 \rightarrow n$, and we miss out two items each time ($\alpha(X_i)\alpha(X_j)$ and $\alpha(X_j)\alpha(X_i)$).
- The second sum is exactly as we saw above – once we combine the two sums, each $\alpha(X_i)$ appears $n-1$ times.
- $\bar{\alpha}$ is a constant, and so simply needs to be multiplied by n .

Overall, we now get:

$$\begin{aligned}
&= \frac{n-1}{n} \left\{ \frac{n-2}{(n-1)^2} \sum_{i \neq j} \alpha(X_i) \alpha(X_j) + \frac{n-1}{(n-1)^2} \sum_{i=1}^n \alpha(X_i)^2 \right. \\
&\quad \left. - \frac{2\bar{\alpha}(n-1)}{n-1} \sum_{i=1}^n \alpha(X_i) + n\bar{\alpha}^2 \right\} \\
&= \frac{n-1}{n} \left\{ \frac{n-2}{(n-1)^2} \sum_{i \neq j} \alpha(X_i) \alpha(X_j) + \frac{1}{n-1} \sum_{i=1}^n \alpha(X_i)^2 \right. \\
&\quad \left. - 2n\bar{\alpha}^2 + n\bar{\alpha}^2 \right\} \\
&= \frac{n-1}{n} \left\{ \frac{n-2}{(n-1)^2} \sum_{i \neq j} \alpha(X_i) \alpha(X_j) + \frac{1}{n-1} \sum_{i=1}^n \alpha(X_i)^2 - n\bar{\alpha}^2 \right\} \\
&= \frac{n-1}{n} \left\{ \frac{n-2}{(n-1)^2} \left(\sum_{i=1}^n \alpha(X_i) \right)^2 + \frac{1}{(n-1)^2} \sum_{i=1}^n \alpha(X_i)^2 - n\bar{\alpha}^2 \right\}
\end{aligned}$$

In the last step, consider the fact that the first term contributes $\frac{n-2}{(n-1)^2}$ lots of $\sum \alpha(x_i)^2$, whereas the second term contributes $\frac{1}{(n-1)^2}$ lots of it. Together, these make

$$\frac{n-2}{(n-1)^2} + \frac{1}{(n-1)^2} = \frac{n-1}{(n-1)^2} = \frac{1}{n-1}$$

Which is indeed what we need, from the second-to-last line. Finally, we write

$$\hat{v}_{\text{jack}} = \frac{1}{n(n-1)} \sum_{i=1}^n (\alpha(X_i) - \bar{\alpha})^2$$

Now, we also have

$$\begin{aligned}
\text{var}(\hat{\theta}) &= \text{var} \left(\mu + \frac{1}{n} \sum_{i=1}^n \alpha(X_i) \right) \\
&= \frac{1}{n^2} \text{var} \left(\sum_{i=1}^n \alpha(X_i) \right)
\end{aligned}$$

Since the X_i are independent, we can write

$$\begin{aligned}
\text{var}(\hat{\theta}) &= \frac{1}{n^2} \sum_{i=1}^n \text{var}(\alpha(X_i)) \\
&= \frac{1}{n^2} \sum_{i=1}^n \mathbb{E} \left(\{ \alpha(X_i) - \bar{\alpha} \}^2 \right)
\end{aligned}$$

Taking the expectation of \hat{v}_{jack} above, we see that it is indeed an unbiased estimator of v in this case. This results only holds true for a linear statistics, but

many smooth statistics can be well approximated by linear statistics. ■

Jackknife estimator of bias: The *jackknife estimator of bias* is given by

$$\widehat{\text{bias}}_{\text{jack}} = (n - 1) \left(\hat{\theta}_{(\text{av})}^{n-1} - \hat{\theta}^n \right)$$

Motivation: We cannot use the same linear statistic as we did in the last proof, because the bias of that statistic is in fact 0. Instead, we need to the *quadratic statistic*

$$\hat{\theta} = \mu + \frac{1}{n} \sum_{i=1}^n \alpha(X_i) + \frac{1}{n^2} \sum_{i=1}^n \sum_{j=1}^n \beta(X_i, X_j)$$

as an estimator of

$$\theta = \mu + \mathbb{E}\{\alpha(X_1)\} + \mathbb{E}\{\beta(X_1, X_2)\} = \mu + a + b$$

We first note that

$$\hat{\theta}_{(-i)}^{n-1} = \mu + \frac{1}{n-1} \sum_{i \neq j} \alpha(X_i) + \frac{1}{(n-1)^2} \sum_{j, k \neq i} \beta(X_i, X_j)$$

Now consider finding the average of each of the terms above

- $\left\langle \frac{1}{n-1} \sum_{i \neq j} \alpha(X_i) \right\rangle = \frac{1}{n} \sum_{i=1}^n \frac{1}{n-1} \sum_{i \neq j} \alpha(X_i)$. Looking at the sums, we find that each $\alpha(X_i)$ is summed over $n - 1$ times. So we can re-write this simply as $\frac{1}{n} \sum_{i=1}^n \alpha(X_i)$.
- $\left\langle \frac{1}{(n-1)^2} \sum_{j, k \neq i} \beta(X_i, X_j) \right\rangle = \frac{1}{n} \sum_{i=1}^n \frac{1}{(n-1)^2} \sum_{j, k \neq i} \beta(X_i, X_j)$

We can split the second sum into two parts:

- Those terms where $i = j$. There are $(n - 1)$ such terms, which gives $\frac{1}{n(n-1)} \sum_{i=j} \beta(X_i, X_j)$.
- Those terms where $i \neq j$. We can simply write those $\sum_{i \neq j} \beta(X_i, X_j)$. The sum above involves $n - 2$ lots of this sum (n from the outer sum, -2 to

exclude each of the subscripts being equal to i).

Thus, all in all, we can write (we omit the arguments of α and β to save space):

$$\hat{\theta}_{(av)}^{n-1} = \mu + \frac{1}{n} \sum_{i \neq j} \alpha + \frac{n-2}{n(n-1)^2} \sum_{i \neq j} \beta + \frac{1}{n(n-1)} \sum_{i=j} \beta$$

It follows that

$$\begin{aligned} \widehat{\text{bias}}_{\text{jack}} &= (n-1)(\hat{\theta}_{(av)}^{n-1} - \hat{\theta}^n) \\ &= (n-1) \left(\frac{n-2}{n(n-1)^2} \sum_{i \neq j} \beta + \frac{1}{n(n-1)} \sum_{i=j} \beta - \frac{1}{n^2} \sum_{i=1}^n \sum_{j=1}^n \beta \right) \end{aligned}$$

The last term, however, can be split as follows

$$\frac{1}{n^2} \sum_{i=1}^n \sum_{j=1}^n \beta = \frac{1}{n^2} \sum_{i=j} \beta + \frac{1}{n^2} \sum_{i \neq j} \beta$$

And so

$$\begin{aligned} \widehat{\text{bias}}_{\text{jack}} &= (n-1)(\hat{\theta}_{(av)}^{n-1} - \hat{\theta}^n) \\ &= (n-1) \left(\left[\frac{n-2}{n(n-1)^2} - \frac{1}{n^2} \right] \sum_{i \neq j} \beta + \left[\frac{1}{n(n-1)} - \frac{1}{n^2} \right] \sum_{i=j} \beta \right) \\ &= (n-1) \left(\left[\frac{n(n-2) - (n-1)^2}{n^2(n-1)^2} \right] \sum_{i \neq j} \beta + \left[\frac{n-(n-1)}{n^2(n-1)} \right] \sum_{i=j} \beta \right) \\ &= -\frac{1}{n^2(n-1)} \sum_{i \neq j} \beta + \frac{1}{n^2} \sum_{i=j} \beta \end{aligned}$$

Taking expectations

$$\begin{aligned} \mathbb{E}(\widehat{\text{bias}}_{\text{jack}}) &= -\frac{1}{n^2(n-1)} n(n-1) \mathbb{E}\{\beta(X_1, X_2)\} \\ &\quad + \frac{1}{n^2} n \mathbb{E}\{\beta(X_1, X_1)\} \\ &= \frac{1}{n} \left[\mathbb{E}\{\beta(X_1, X_1)\} - \mathbb{E}\{\beta(X_1, X_2)\} \right] \\ &= \frac{1}{n} \left[\mathbb{E}\{\beta(X_1, X_1)\} - b \right] \\ &= \text{bias}(\hat{\theta}_n) \end{aligned}$$

As expected. ■

The jackknife can also be used to generate $B = n$ “samples” from X_1, \dots, X_n by missing out a single element from each sample. However, if n is small, B will be small, and in the context of a Monte Carlo test, this will severely limit the possible choices of significance level α .

The Bootstrap

The non-parametric bootstrap is another form of resampling which involves forming samples *with* replacement from observed samples.

Consider a set of independent random variables $X = (X_1, \dots, X_n)$, with distribution function F , of which we have a realisation x_1, \dots, x_n . We are interested in the distribution of a root or pivot $R_n(X; F)$, which we denote $K_n(F)$. Such a “root” could include any statistic $T_n(X)$, but also quantities like $\frac{\sqrt{n}}{\hat{\sigma}} \{ \bar{X} - \mathbb{E}_F(X_1) \}$.

Our strategy will be to use a different distribution function \hat{F} that approximates F , and then estimate $K_n(F)$ by $K_n(\hat{F})$, as follows

Bootstrap estimator: We estimate $K_n(F)$ by $K_n(\hat{F})$, which we find as follows

1. Draw B independent *bootstrap samples*, each of size n

$$X_b^* = (X_{b1}^*, \dots, X_{bn}^*) \quad b = 1, \dots, B$$

where each of the X^* are independently drawn from the distribution \hat{F} .

2. Approximate $K_n(\hat{F})$ by the ECDF of

$$\{R_n(X_b^*; \hat{F}) : b = 1, \dots, B\}$$

In general, $B = 100$ or $B = 200$ is often sufficient to estimate a variance or quantile, but $B = 1000$ is recommended to estimate the entire distribution.

Notes: The question of how to choose \hat{F} remains.

- In a parametric model $\mathcal{F} = \{F_\theta : \theta \in \Theta\}$, the *parametric bootstrap* uses $\hat{F} = F_{\hat{\theta}}$, where $\hat{\theta}$ is the MLE of θ .
- The *non-parametric bootstrap* simply uses the ECDF \hat{F}_n of $X = (X_1, \dots, X_n)$. In that case, step 1 above simply involves taking a re-sample of size n with replacement from

x_1, \dots, x_n , with equal probability placed on each expression.

Notes: If N_i denotes the number of times x_i appears in the bootstrap sample, then the vector (N_1, \dots, N_n) (note that $N_+ = 1$) has a one-to-one mapping with the $(2n - 1)$ binary tuple that contains N_1 ones, followed by a zero, followed by N_2 ones, etc... There are ${}^{2n-1}C_n$ such tuples, and this is therefore the number of bootstrap samples.

Furthermore, $N_i \sim \text{Multi}\left(n; \frac{1}{n}, \dots, \frac{1}{n}\right)$, and so

$$\mathbb{P}(N_1 = n_1, \dots, N_n = n_n) = \frac{n!}{n_1! \dots n_n!} \left(\frac{1}{n}\right)^n \leq \frac{n!}{n^n}$$

So the most likely sample is the original data, with probability $n!/n^n$.

R-CODE: The bootstrap samples themselves can be taken using

```
X.star <- matrix(NA, nrow=B, ncol=n)
For(b in 1:B) {
  X.star[b,] <- sample(x, n, replace=TRUE)
}
```

Where \mathbf{x} contains our data, n is the number of data items, and B is the number of bootstrap samples we want to take.

We cover a few examples to clarify this concept:

1. Finding a confidence interval for $\mu = \mathbb{E}_F\{X_1\}$

An analytic approach would be to define the root

$$R_n(X; F) = \sqrt{n}(\bar{X} - \mu) \rightarrow^d N(0, \text{Var}_F(X_1))$$

and then to estimate⁴

⁴ To prove this result, consider that the ECDF of X_1^* , and associated density, are

$$\hat{F}_n(x) = \frac{1}{n} \sum_{i=1}^n \mathbb{I}_{\{x_i \leq x\}} \quad \hat{f}_n(x) = \frac{1}{n} \sum_{i=1}^n \mathbb{I}_{\{x_i = x\}}$$

and so

$$\mathbb{E}(X_1^*) = \int x \hat{f}_n(x) dx = \int \frac{1}{n} \sum_{i=1}^n x \mathbb{I}_{\{x_i = x\}} dx = \frac{1}{n} \sum_{i=1}^n x_i = \bar{x}$$

and similarly

$$\text{Var}_F(X_1) = \text{Var}_{\hat{F}}(X_1^*) = \frac{1}{n} \sum (x_i - \bar{x})^2$$

and finally estimating the $(1 - \alpha)$ confidence interval as

$$\left[\bar{x} - \frac{z_{\sigma^2, \frac{1}{2}\alpha}}{\sqrt{n}}, \bar{x} + \frac{z_{\sigma^2, \frac{1}{2}\alpha}}{\sqrt{n}} \right]$$

Where $z_{\sigma^2, \alpha}$ **denotes the upper** α point of a $N(0, \sigma^2)$ distribution.

In reality, however, we do *not* know the distribution of the root. We therefore use the following bootstrap method:

- Fix a large B such that $\frac{1}{2}\alpha(B + 1)$ is an integer.
- For each, b , generate independent bootstrap samples X_b^* by re-sampling, and compute

$$R_b^* = R_n(X_b^*; \hat{F}_n) = \sqrt{n}(\bar{X}_b^* - \bar{X})$$

- Approximate the $(1 - \alpha)$ confidence interval as

$$\left[\bar{x} - R_{(\lfloor \frac{1-\alpha}{2} \rfloor [B+1])}^*, \bar{x} + R_{(\lceil \frac{\alpha}{2} \rceil [B+1])}^* \right]$$

More generally, the bootstrap confidence intervals for a parameter $N(0, \sigma^2)$ are often based on the root $R_n(X; F) = \sqrt{n}(\tilde{\theta} - \theta)$, where $\tilde{\theta}$ is an estimator of θ . These are called *percentile intervals*.

2. Estimating $\text{Var}_F\{\tilde{\theta}(X)\}$, where $\tilde{\theta}(X)$ is an estimator of $\theta = \theta(F)$

Our step should be to try and calculate the non-parametric bootstrap estimator $\text{Var}_{\hat{F}_n}\{\tilde{\theta}(X^*)\}$ analytically, then we apply the following algorithm

- Generate B independent bootstrap samples X_b^* (by re-sampling)
- For each bootstrap sample, calculate $\tilde{\theta}_b^* = \tilde{\theta}(X_b^*)$
- Approximate $\text{Var}_{\hat{F}_n}\{\tilde{\theta}(X^*)\}$ by

$$\frac{1}{B-1} \sum_{b=1}^B (\tilde{\theta}_b^* - \bar{\tilde{\theta}}^*)^2 \quad \text{where } \bar{\tilde{\theta}}^* = \frac{1}{B} \sum_{i=1}^B \tilde{\theta}_i^*$$

$$\text{Var}(X_1^*) = \int (x - \bar{x})^2 \hat{f}_n(x) dx = \int \frac{1}{n} \sum_{i=1}^n (x - \bar{x})^2 \mathbb{I}_{\{x_i=x\}} dx = \frac{1}{n} \sum_{i=1}^n (x_i - \bar{x})^2$$

3. Estimating the bias of $\tilde{\theta}(X)$, an estimator of $\theta = \theta(F)$

The bias of $\tilde{\theta}$ is

$$\text{bias}_F \{ \tilde{\theta}(X) \} = \mathbb{E}_F \{ \tilde{\theta}(X) \} - \theta$$

the non-parametric bootstrap estimator of the bias is

$$\text{bias}_{\hat{F}_n} \{ \tilde{\theta}(X^*) \} = \mathbb{E}_{\hat{F}_n} \{ \tilde{\theta}(X^*) \} - \hat{\theta} \quad \text{where } \hat{\theta} = \theta(\hat{F}_n)$$

If we cannot compute this expression directly, we can apply the following algorithm instead:

- Generate B independent bootstrap samples X_b^* (by re-sampling)
- For each bootstrap sample, calculate $\tilde{\theta}_b^* = \tilde{\theta}(X_b^*)$
- Approximate $\text{bias}_{\hat{F}_n}(\tilde{\theta})$ by

$$\frac{1}{B} \sum_{b=1}^B (\tilde{\theta}_b^* - \hat{\theta})$$

4. Estimating the distribution of the sample median, $F^{-1}(\frac{1}{2})$

In this case, the distribution can be calculated analytically. Let $x_{(j)}$ be the j^{th} ordered statistic of x_1, \dots, x_n , the *original* sample. Let also \hat{F}_n^* be the empirical distribution function of a given bootstrap sample. Then

$$\begin{aligned} \mathbb{P} \left(\hat{F}_n^{*-1}(\frac{1}{2}) \leq x_{(j)} \right) &= \mathbb{P} \left(\sum_{i=1}^n \mathbb{I}_{\{X_i^* \leq x_{(j)}\}} \geq \lfloor \frac{1}{2} n \rfloor \right) \\ &= \sum_{k=\lfloor \frac{1}{2} n \rfloor}^n \binom{n}{k} \left(\frac{j}{n} \right)^k \left(\frac{n-j}{n} \right)^{n-k} \end{aligned}$$

Effectively, this simply the probability of our bootstrap sample containing more than $n/2$ items which are $\leq x_{(j)}$ [each with probability j/n]. The sum above is over all the possible number of “more than $n/2$ ” items. From this, we have, for $j = 2, \dots, n$

$$\mathbb{P} \left(\hat{F}_n^{*-1}(\frac{1}{2}) = x_{(j)} \right) = \mathbb{P} \left(\hat{F}_n^{*-1}(\frac{1}{2}) \leq x_{(j)} \right) - \mathbb{P} \left(\hat{F}_n^{*-1}(\frac{1}{2}) \leq x_{(j-1)} \right)$$

and

$$\mathbb{P} \left(\hat{F}_n^{*-1}(\frac{1}{2}) = x_{(1)} \right) = \sum_{k=\lfloor \frac{1}{2} n \rfloor}^n \binom{n}{k} \left(\frac{1}{n} \right)^k \left(\frac{n-1}{n} \right)^{n-k}$$

We now consider a final interesting example in which the nonparametric bootstrap *fails*. Consider $X_1, \dots, X_n \sim U(0, \theta]$, and define the root

$$R_n(X; F_\theta) = \frac{n(\theta - X_{(n)})}{\theta}$$

where $X_n = \max(X_1, \dots, X_n)$.

We first find the distribution of $R_n(X; F)$

$$\mathbb{P}_F(R_n > x) = \mathbb{P}_F\left(X_{(n)} < \frac{(n-x)\theta}{n}\right) = \left(1 - \frac{x}{n}\right)^n \rightarrow e^{-x}$$

So $R_n \sim \text{Exp}(1)$. Note, however, that if we use \hat{F}_n as our distribution, then our best estimator of $\theta = x_{(n)}$, and so

$$R_n(X^*, \hat{F}_n) = \frac{n(x_{(n)} - X_{(n)}^*)}{x_{(n)}}$$

And so

$$\begin{aligned} \mathbb{P}_{\hat{F}_n}(R_n = 0) &= \mathbb{P}(X_{(n)}^* = x_{(n)}) \\ &= \mathbb{P}(x_{(n)} \in \text{bootstrap sample}) \\ &= 1 - \mathbb{P}(x_{(n)} \notin \text{bootstrap sample}) \\ &= 1 - \left[\mathbb{P}(x_{(n)} \neq i^{\text{th}} \text{ element in sample})\right]^n \\ &= 1 - \left(1 - \frac{1}{n}\right)^n \\ &\rightarrow 1 - e^{-1} \end{aligned}$$

Clearly, therefore, the asymptotic distribution of the nonparametric bootstrap here does not tend to that of the root.

Let us consider, instead, the parametric bootstrap, which uses the MLE of $\hat{\theta} = x_{(n)}$:

$$R_n(X^*, F_{\hat{\theta}}) = \frac{n(\hat{\theta} - X_{(n)}^*)}{\hat{\theta}}$$

Now

$$\mathbb{P}_{F_{\hat{\theta}}}(R_n > x) = \mathbb{P}_{F_{\hat{\theta}}}\left(X_{(n)}^* < \frac{(n-x)\hat{\theta}}{n}\right) = \left(1 - \frac{x}{n}\right)^n \rightarrow e^{-x}$$

Which does indeed tend to the distribution of the root.

Note the key difference between \hat{F}_n , which picks each value with equal probability, and F_θ , in which the probability of picking values in an interval depends on the value itself.

This behaviour is due to the non-standard asymptotics of the uniform distribution.

Bayesian Inference

Bayes' Theorem states that

$$\pi(\boldsymbol{\theta} | \mathbf{x}) = \frac{L(\mathbf{x} | \boldsymbol{\theta})p(\boldsymbol{\theta})}{\int L(\mathbf{x} | \boldsymbol{\theta})p(\boldsymbol{\theta}) d\boldsymbol{\theta}}$$

Where

- $\pi(\boldsymbol{\theta} | \mathbf{x})$ is the *posterior distribution* of $\boldsymbol{\theta}$
- $L(\mathbf{x} | \boldsymbol{\theta})$ is the *likelihood*
- $p(\boldsymbol{\theta})$ is the *prior* on $\boldsymbol{\theta}$.

EXAMPLE: Consider an auto-regressive (times-series) model of order k , in which we observe data $\mathbf{x} = (x_1, \dots, x_N)$, which, for $t > k + 1$ is generated by the process

$$x_t = \sum_{r=1}^k a_r x_{t-r} + \varepsilon_t \quad \varepsilon_t \sim N(0, \sigma^2)$$

We can express this as

$$\mathbf{x}_k = X_k \mathbf{a} + \boldsymbol{\varepsilon}$$

Where

- $\mathbf{x}_k = (x_{k+1}, \dots, x_N)^T$
- X_k is an $(N - k) \times k$ matrix with $(X_k)_{ij} = x_{k+i-j}$.
- $\mathbf{a} = (a_1, \dots, a_k)$
- $\boldsymbol{\varepsilon} = (\varepsilon_{k+1}, \dots, \varepsilon_N)^T$

We then have

$$L(\mathbf{x}_k | \mathbf{a}, \sigma^2) = \frac{1}{(2\pi\sigma^2)^{\frac{1}{2}(n-k)}} \exp\left(-\frac{1}{2}(\mathbf{x}_k - X_k \mathbf{a})^T (\sigma^2 I)^{-1} (\mathbf{x}_k - X_k \mathbf{a})\right)$$

We might place the following priors on \mathbf{a} and σ^2

$$\mathbf{a} \sim N_k(\boldsymbol{\mu}_a, \Sigma_a)$$

$$\Rightarrow p(\mathbf{a}) = \frac{1}{(2\pi)^{k/2} \sqrt{|\Sigma_a|}} \exp\left(-\frac{1}{2}(\mathbf{a} - \boldsymbol{\mu}_a)^T \Sigma_a^{-1} (\mathbf{a} - \boldsymbol{\mu}_a)\right)$$

$$\sigma^2 \sim \text{InverseGamma}(\alpha, \beta)$$

$$\Rightarrow p(\sigma^2) = \frac{\beta^\alpha}{\Gamma(\alpha)} (\sigma^2)^{-\alpha-1} e^{-\beta/\sigma^2}$$

We are interested here in calculating marginal summary statistics for θ . For example, the posterior expectation

$$\mathbb{E}_\pi \{ \phi(\theta) \} = \int \phi(\theta) \pi(\theta | \mathbf{x}) \, d\theta$$

Recall that, using Monte-Carlo integration, we can approximate

$$\mathbb{E}_\pi \{ \phi(\theta) \} \approx \frac{1}{n} \sum_{i=1}^n \phi(\theta_i) \quad \theta_1, \dots, \theta_n \sim \pi(\theta | \mathbf{x})$$

Unfortunately π may be difficult to sample from, since it is often high dimensional and of unfamiliar form. We resolve this problem as follows:

- We construct a Markov chain $\theta^{(0)}, \theta^{(1)}, \dots$ which has π as its stationary distribution. In other words, such that

$$\left\| \mathcal{K}^t(\theta^{(0)}, ?) - \pi(? | \mathbf{x}) \right\| \rightarrow 0 \quad \text{as } t \rightarrow \infty$$

Where $\mathcal{K}^t(\theta^{(0)}, ?) = \mathbb{P}(\theta^{(t)} = ? | \theta^{(0)})$.

- We run the chain until it reaches equilibrium.
- Further realisations can be regarded as a dependent sample from π .

This method is called **Markov Chain Monte Carlo**. We consider various methods to construct such a chain.

The Gibbs Sampler

Gibbs Sampling: Imagine the random vector we are sampling has p dimensions, and distribution $\pi(\theta)$. The Gibbs sampler samples from this distribution as follows:

1. Begin in some arbitrary state $\theta^{(0)} = (\theta_1^{(0)}, \dots, \theta_p^{(0)})$.

2. At time t , in state $\theta^{(t)}$, update the state vector one components at a time:

$$\begin{aligned} \theta_1^{(t+1)} &\sim \pi(\theta_1 | \theta_1^{(t)}, \theta_2^{(t)}, \dots, \theta_p^{(t)}) \\ \theta_2^{(t+1)} &\sim \pi(\theta_2 | \theta_1^{(t+1)}, \theta_2^{(t)}, \theta_3^{(t)}, \dots, \theta_p^{(t)}) \\ \theta_p^{(t+1)} &\sim \pi(\theta_p | \theta_1^{(t+1)}, \dots, \theta_{p-1}^{(t+1)}, \theta_p^{(t)}) \end{aligned}$$

3. Collect a total of T samples.
4. Discard the first b samples as “*burn-in*”.
5. Treat $\left\{ \theta^{(t)} \right\}_{t=b+1}^T$ as a dependent sample from π .

The distribution of one component conditional on another is called the **full joint conditional distribution**.

EXAMPLE: Consider the AR model from above. The full joint conditions for \mathbf{a} and σ^2 are

$$\begin{aligned}\pi(\sigma^2 \mid \mathbf{x}_k, \mathbf{a}) &\propto L(\mathbf{x}_k \mid \mathbf{a}, \sigma^2) \pi(\sigma^2) \\ &\propto \frac{1}{(\sigma^2)^{\frac{1}{2}(N-k)}} \exp\left(-\frac{1}{2\sigma^2} (\mathbf{x}_k - X_k \mathbf{a})^T (\mathbf{x}_k - X_k \mathbf{a})\right) \\ &\quad (\sigma^2)^{-\alpha-1} e^{-\beta/\sigma^2} \\ &\propto \frac{1}{(\sigma^2)^{\frac{1}{2}(N-k)}} (\sigma^2)^{-\alpha-\frac{1}{2}(N-k)-1} \\ &\quad \exp\left(-\frac{1}{\sigma^2} \left\{ \frac{1}{2} (\mathbf{x}_k - X_k \mathbf{a})^T (\mathbf{x}_k - X_k \mathbf{a}) + \beta \right\}\right) \\ &\sim \text{InverseGamma} \left(\begin{array}{l} \alpha + \frac{1}{2}(N-k), \\ \beta + \frac{1}{2} (\mathbf{x}_k - X_k \mathbf{a})^T (\mathbf{x}_k - X_k \mathbf{a}) \end{array} \right)\end{aligned}$$

And

$$\begin{aligned}\pi(\mathbf{a} \mid \mathbf{x}_k, \sigma^2) &= L(\mathbf{x}_k \mid \mathbf{a}, \sigma^2) \pi(\mathbf{a}) \\ &\propto \exp\left(-\frac{1}{2\sigma^2} (\mathbf{x}_k - X_k \mathbf{a})^T (\mathbf{x}_k - X_k \mathbf{a})\right) \\ &\quad \exp\left(-\frac{1}{2} (\mathbf{a} - \boldsymbol{\mu}_a)^T \Sigma_a^{-1} (\mathbf{a} - \boldsymbol{\mu}_a)\right)\end{aligned}$$

Since we want a quadratic in \mathbf{a} , we write

$$\begin{aligned}\pi(\mathbf{a} \mid \mathbf{x}_k, \sigma^2) &\propto \exp\left(-\frac{1}{2} \left[\mathbf{a}^T \left(-\sigma^{-2} X_k^T X_k + \Sigma_a^{-1} \right) \mathbf{a} \right. \right. \\ &\quad \left. \left. - 2\mathbf{a}^T \left(\sigma^{-2} X_k^T \mathbf{x}_k + \Sigma_a^{-1} \boldsymbol{\mu}_a \right) \right] \right) \\ &\propto \exp\left(-\frac{1}{2} \left\{ \mathbf{a}^T \Sigma^{-1} \mathbf{a} - 2\mathbf{a}^T \Sigma^{-1} \boldsymbol{\mu} \right\}\right) \\ &\propto \exp\left(-\frac{1}{2} (\mathbf{a} - \boldsymbol{\mu})^T \Sigma^{-1} (\mathbf{a} - \boldsymbol{\mu})\right) \\ &\sim N_k(\boldsymbol{\mu}, \Sigma)\end{aligned}$$

With

$$\Sigma^{-1} = \sigma^{-2} X_k^T X_k + \Sigma_a^{-1} \quad \boldsymbol{\mu} = \Sigma(\sigma^{-2} X_k^T \mathbf{x}_k + \Sigma_a^{-1} \boldsymbol{\mu}_a)$$

The Gibbs sampler can be used to great effect when

- There is some missing data

- The likelihood is difficult to use in its native form, but easy to use when conditional on some unobserved data (which we can treat as “missing” data).

In both the cases above, we can explicitly calculate $L(y_{\text{obs}}, y_{\text{miss}} | \theta)$. We then simulate θ and y_{miss} together using a Gibbs sampler, with

$$\pi(\theta, y_{\text{miss}} | y_{\text{obs}}) \propto L(y_{\text{obs}}, y_{\text{miss}} | \theta) p(\theta)$$

Of course, in reality, the distribution we are interested in is

$$\pi(\theta | y_{\text{obs}}) = \int \pi(\theta, y_{\text{miss}} | y_{\text{obs}}) dy_{\text{miss}}$$

The values of θ returned by the Gibbs sampler, however, have precisely that distribution, because by ignoring the values of y_{miss} that the Gibbs sampler returns, we are effectively “integrating over” y_{miss} .

EXAMPLE: Consider a group of N animals, each assigned to one of four categories

$$(y_1, y_2, y_3, y_4)$$

with probabilities

$$\left(\frac{2+\theta}{4}, \frac{1-\theta}{4}, \frac{1-\theta}{4}, \frac{\theta}{4}\right)$$

The likelihood is then multinomial

$$L(\mathbf{y} | \theta) \propto \left(\frac{2+\theta}{4}\right)^{y_1} \left(\frac{1-\theta}{4}\right)^{y_2+y_3} \left(\frac{\theta}{4}\right)^{y_4}$$

Suppose we place a $\text{Beta}(\alpha, \beta)$ prior on θ . Then

$$\begin{aligned} \pi(\theta | \mathbf{y}) &\propto L(\mathbf{y} | \theta) p(\theta) \\ &\propto \left(\frac{2+\theta}{4}\right)^{y_1} \left(\frac{1-\theta}{4}\right)^{y_2+y_3} \left(\frac{\theta}{4}\right)^{y_4} \theta^{\alpha-1} (1-\theta)^{\beta-1} \\ &\propto (2+\theta)^{y_1} (1-\theta)^{y_2+y_3+\beta-1} \theta^{y_4+\alpha-1} \end{aligned}$$

This is hard to sample from.

Consider, instead, splitting our data into the following *five* groups

$$(y_1 - z, z, y_2, y_3, y_4)$$

With probabilities

$$\left(\frac{1}{2}, \frac{\theta}{4}, \frac{1-\theta}{4}, \frac{1-\theta}{4}, \frac{\theta}{4}\right)$$

We then have

$$\begin{aligned} \pi(\theta, z | \mathbf{y}) &\propto L(\mathbf{y} | \theta, z) p(\theta) \\ &\propto \frac{N!}{z!(y_1-z)! \dots} \left(\frac{1}{2}\right)^{y_1-z} \left(\frac{\theta}{4}\right)^z \left(\frac{1-\theta}{4}\right)^{y_2+y_3} \left(\frac{\theta}{4}\right)^{y_4} \theta^{\alpha-1} (1-\theta)^{\beta-1} \\ &\propto {}^{y_1}C_z \left(\frac{1}{2}\right)^{y_1-z} \left(\frac{\theta}{4}\right)^z \theta^{y_4+\alpha-1} (1-\theta)^{y_2+y_3+\beta-1} \end{aligned}$$

[Note – it is important, in this case, to keep the factor of $\theta / 4$ intact, because getting rid of the $1/4$ would imply getting rid of a factor of $1/4^z$. Since, however, we will be needing the distribution of z , this is not legitimate].

And so

$$\theta | z, \mathbf{y} \sim \text{Beta}(z + y_4 + \alpha, y_2 + y_3 + \beta)$$

Finding the joint conditional for z is slightly harder. Here are two ways to do it

1. Note that

$$\pi(\theta, z | \mathbf{y}) \propto {}^{y_1}C_z \left(\frac{1}{2}\right)^{y_1-z} \left(\frac{\theta}{4}\right)^z$$

This looks like a binomial, but $1/2$ and $\theta / 4$ don't sum to 1. We'll therefore try and find something we can multiply both of them by to get something that *does* sum to 1.

$$\begin{aligned} \frac{1}{2} ? + \frac{\theta}{4} ? &= 1 \\ ? &= \frac{1}{\frac{\theta}{4} + \frac{1}{2}} = \frac{4}{2 + \theta} \end{aligned}$$

And so

$$\pi(\theta, z | \mathbf{y}) \propto {}^{y_1}C_z \left(\frac{2}{2+\theta}\right)^{y_1-z} \left(\frac{\theta}{2+\theta}\right)^z$$

As such

$$z | \theta, \mathbf{y} \sim \text{Bin}\left(y_1, \frac{\theta}{2+\theta}\right)$$

2. Or, we can simply note that once θ is known, the conditional probability of an observation being in z given it's in y_1 is

$$\frac{\mathbb{P}(\in z)}{\mathbb{P}(\in y_1)} = \frac{\frac{\theta}{4}}{\frac{1}{2} + \frac{\theta}{4}} = \frac{\theta}{2 + \theta}$$

Which yields the same result.

The Metropolis-Hastings Algorithm

The Gibbs sampler is a special case of the Metropolis-Hastings (MH) algorithm. The algorithm basically samples values from an “approximate” distribution and then “corrects” these so that they asymptotically behave as if they came from the stationary distribution.

The Metropolis Hastings (MH) algorithm:

1. Begin with some arbitrary state $\theta^{(0)}$.
2. Simulate $\phi \sim q(\theta^{(t)}, \phi)$; q is our “approximate transition probability” from $\theta^{(t)}$ to ϕ .
3. Let

$$\alpha(\theta^{(t)}, \phi) = \min\left(1, \frac{\pi(\phi | \mathbf{x})q(\phi, \theta^{(t)})}{\pi(\theta^{(t)} | \mathbf{x})q(\theta^{(t)}, \phi)}\right)$$

4. Set $\theta^{(t+1)} = \phi$ with probability $\alpha(\theta^{(t)}, \phi)$, or else reject it and set $\theta^{(t+1)} = \theta^{(t)}$.

So the Markov transition Kernel for the chain is given by

$$\mathcal{P}_H(\theta, B) = \underbrace{\int_B \mathcal{K}_H(\theta, \phi) d\phi}_{\text{Moving forward to some } \phi \in B} + \underbrace{r(\theta)I_B(\theta)}_{\text{Staying, if it so happened that } \theta \text{ was already } \in B} + \underbrace{0}_{\text{We neither move to any state in } B, \text{ nor do we remain there}}$$

Where

$$\begin{aligned} \mathcal{K}_H(\theta, \phi) &= q(\theta, \phi)\alpha(\theta, \phi) \\ r_B(\theta) &= 1 - \int_{\forall x} \mathcal{K}_H(\theta, x) dx \end{aligned}$$

This leads to values that behave as if they come from $\pi(\theta)$.

Proof: First, note that

$$\begin{aligned} \pi(x)\mathcal{K}_H(x, y) &= \pi(x)q(x, y) \min\left(1, \frac{\pi(y)q(y, x)}{\pi(x)q(x, y)}\right) \\ &= \min\left(\pi(x)q(x, y), \pi(y)q(y, x)\right) \\ &= \min\left(\frac{\pi(x)q(x, y)}{\pi(y)q(y, x)}, 1\right)\pi(y)q(y, x) \\ &= \pi(y)\mathcal{K}_H(y, x) \end{aligned}$$

This implies that

$$\begin{aligned} \int_{\forall x} \pi(x)\mathcal{K}_H(x, y) dx &= \pi(y) \int_{\forall x} \mathcal{K}_H(y, x) dx \\ &= \pi(y)[1 - r(y)] \end{aligned}$$

With this key result in hand, we are now ready to prove our theorem. Effectively, we want to show that if we sample \mathbf{x} from the stationary distribution, the probability of the Markov chain moving us from \mathbf{x} to any element in B is the same as the probability of being in B under the posterior distribution. In other words, we want to show that

$$\int_{\forall x} \pi(x) \mathcal{P}_H(x, B) \, dx = \int_B \pi(y) \, dy$$

To do that, consider

$$\begin{aligned} & \int_{\forall x} \pi(x) \mathcal{P}_H(x, B) \, dx \\ &= \int_{\forall x} \pi(x) \left[\int_B \mathcal{K}_H(x, y) \, dy + r(x) I_B(x) \right] \, dx \\ &= \int_{\forall x} \int_B \pi(x) \mathcal{K}_H(x, y) \, dy \, dx + \int_{\forall x} \pi(x) r(x) I_B(x) \, dx \\ &= \int_B \int_{\forall x} \pi(x) \mathcal{K}_H(x, y) \, dx \, dy + \int_B \pi(x) r(x) \, dx \end{aligned}$$

We can now use our previously derived result for $\int_{\forall x} \pi(x) \mathcal{K}_H(x, y)$, which gives

$$\begin{aligned} & \int_{\forall x} \pi(x) \mathcal{P}_H(x, B) \, dx \\ &= \int_B \pi(y) [1 - r(y)] \, dy + \int_B \pi(x) r(x) \, dx \\ &= \int_B \pi(y) \, dy \end{aligned}$$

As required. ■

Note: A great advantage of the MH algorithm over Gibbs sampling is that we do not need to know the normalisation constant of π . Also, we do not need to know any of the joint conditionals.

The acceptance function might be easier to understand when written as follows

$$\alpha(\boldsymbol{\theta}^{(t)}, \boldsymbol{\phi}) = \min \left(1, \frac{\pi(\boldsymbol{\phi} | \mathbf{x})}{\pi(\boldsymbol{\theta}^{(t)} | \mathbf{x})} \div \frac{q(\boldsymbol{\theta}^{(t)}, \boldsymbol{\phi})}{q(\boldsymbol{\phi}, \boldsymbol{\theta}^{(t)})} \right)$$

This makes it clearer that it is indeed in the form (actual density / proposal density) that we saw in importance sampling.

Note that when the proposal distribution q is symmetric, it is called a *symmetric* (or *Metropolis*) proposal. In that case, $q(\boldsymbol{\phi}, \boldsymbol{\theta}) = q(\boldsymbol{\theta}, \boldsymbol{\phi})$ and so the acceptance probability reduces to

$$\alpha(\boldsymbol{\theta}, \boldsymbol{\phi}) = \min \left\{ 1, \frac{\pi(\boldsymbol{\phi} | \boldsymbol{x})}{\pi(\boldsymbol{\theta} | \boldsymbol{x})} \right\}$$

There are a number of possible choices for the proposal distribution q . Here are two common examples:

- **Random-Walk (RW) Metropolis**

Here, we specify

$$\boldsymbol{\phi} = \boldsymbol{\theta} + \boldsymbol{z} \quad \boldsymbol{z} \sim f$$

Common choices for f may include a uniform distribution or a multivariate normal. The distribution is often – but not always – chosen to be symmetric. In RW metropolis, we are effectively choosing our new value to be close to the current one, the reasoning being that when the stationary distribution has been reached, we're more likely to spend longer at points of high density.

- **Independence Sampler**

Here, the candidate observation is drawn independently of the current state, so that $q(\boldsymbol{\theta}, \boldsymbol{\phi}) = g(\boldsymbol{\phi})$. The corresponding acceptance probability can be written

$$\alpha(\boldsymbol{\theta}, \boldsymbol{\phi}) = \min \left\{ 1, \frac{\omega(\boldsymbol{\phi})}{\omega(\boldsymbol{\theta})} \right\} \quad \omega(\cdot) = \frac{\pi(\cdot | \boldsymbol{x})}{g(\cdot)}$$

This is precisely the importance weight function that would be used in importance sampling, given observations from g being used to sample from π .

Theorem: The Gibbs sampler is a special case of the Metropolis Hastings algorithm.

Proof: Suppose we have a current estimate of our parameters $\boldsymbol{\theta} = (\theta_1, \dots, \theta_p)$. We break each iteration of the MH algorithm into steps, each of which update a single value of $\boldsymbol{\theta}$.

1. Start by setting $\boldsymbol{\theta}^{(0)} = \boldsymbol{\theta}$. Set $t = 1$.

2. Propose a new value ϕ for a single component θ_t , so produce a new vector

$$\boldsymbol{\theta}^{(t)} = \left(\theta_1^{(t-1)}, \dots, \theta_{t-1}^{(t-1)}, \phi, \theta_{t+1}^{(t-1)}, \dots, \theta_p^{(t-1)} \right)$$

with proposal density

$$q_t(\boldsymbol{\theta}^{(t-1)}, \boldsymbol{\theta}^{(t)}) = \pi(\phi \mid \boldsymbol{\theta}_{(-t)}^{(t-1)}, \mathbf{x}) = \pi(\theta_t^{(t)} \mid \boldsymbol{\theta}_{(-t)}^{(t-1)}, \mathbf{x})$$

Where $\boldsymbol{\theta}_{(-t)}^{(t-1)} = \left(\theta_1^{(t-1)}, \dots, \cancel{\theta_t^{(t-1)}}, \dots, \theta_p^{(t-1)} \right)$. Denote

3. Accept the proposal with acceptance probability $\alpha(\boldsymbol{\theta}^{(t-1)}, \boldsymbol{\theta}^{(t)})$.
 - a. If it is accepted, keep $\boldsymbol{\theta}^{(t)}$.
 - b. Otherwise, keep the previous iteration $\boldsymbol{\theta}^{(t)} \leftarrow \boldsymbol{\theta}^{(t-1)}$.

All we need to do is to show that the acceptance probability in this case is always 1. In that case, the algorithm above is equivalent to the Gibbs sampler.

Now, $\alpha(\boldsymbol{\theta}^{(t-1)}, \boldsymbol{\theta}^{(t)}) = \min\{1, A_t\}$, where

$$\begin{aligned} A_t &= \frac{\pi(\boldsymbol{\theta}^{(t)} \mid \mathbf{x}) q_t(\boldsymbol{\theta}^{(t)}, \boldsymbol{\theta}^{(t-1)})}{\pi(\boldsymbol{\theta}^{(t-1)} \mid \mathbf{x}) q_t(\boldsymbol{\theta}^{(t-1)}, \boldsymbol{\theta}^{(t)})} \\ &= \frac{\pi(\boldsymbol{\theta}^{(t)} \mid \mathbf{x}) \pi(\theta_t^{(t-1)} \mid \boldsymbol{\theta}_{(-t)}^{(t)}, \mathbf{x})}{\pi(\boldsymbol{\theta}^{(t-1)} \mid \mathbf{x}) \pi(\theta_t^{(t)} \mid \boldsymbol{\theta}_{(-t)}^{(t-1)}, \mathbf{x})} \\ &= \frac{\pi(\boldsymbol{\theta}^{(t)} \mid \mathbf{x}) / \pi(\theta_t^{(t)} \mid \boldsymbol{\theta}_{(-t)}^{(t-1)}, \mathbf{x})}{\pi(\boldsymbol{\theta}^{(t-1)} \mid \mathbf{x}) / \pi(\theta_t^{(t-1)} \mid \boldsymbol{\theta}_{(-t)}^{(t)}, \mathbf{x})} \end{aligned}$$

Note, however, that $\boldsymbol{\theta}_{(-t)}^{(t)} = \boldsymbol{\theta}_{(-t)}^{(t-1)}$, so

$$\begin{aligned} A_t &= \frac{\pi(\boldsymbol{\theta}^{(t)} \mid \mathbf{x}) / \pi(\theta_t^{(t)} \mid \boldsymbol{\theta}_{(-t)}^{(t)}, \mathbf{x})}{\pi(\boldsymbol{\theta}^{(t-1)} \mid \mathbf{x}) / \pi(\theta_t^{(t-1)} \mid \boldsymbol{\theta}_{(-t)}^{(t-1)}, \mathbf{x})} \\ &= \frac{\pi(\boldsymbol{\theta}_{(-t)}^{(t)})}{\pi(\boldsymbol{\theta}_{(-t)}^{(t-1)})} \\ &= 1 \end{aligned}$$

(To go from the first line to the second line, use Bayes' Theorem on conditional π probabilities). As required. ■

It is generally accepted that if one wants T independent samples from π , but must instead resort to obtaining T depend samples (eg: via MCMC), then

somehow, the resulting effective sample size (relative to independent sampling) is less than T , due to the autocorrelation of the chain.

Radford Neal defined

Effective sample size due to autocorrelation: (for marginally scalar samples)

$$\text{ESS}(\boldsymbol{\theta}) = \frac{T}{1 + 2\sum_{\ell=1}^{T-1} \hat{\rho}(\ell)}$$

Where $\hat{\rho}(\ell)$ is the sample autocorrelation at lag ℓ :

$$\hat{\rho}(\ell) = \frac{\hat{\gamma}(\ell)}{\hat{\gamma}(0)}$$

where

$$\hat{\gamma}(\ell) = \frac{1}{T-\ell} \sum_{t=1}^{T-\ell} (\theta^{(t)} - \bar{\theta})(\theta^{(t+\ell)} - \bar{\theta})$$

$$\bar{\theta} = \frac{1}{T} \sum_{t=1}^T \theta^{(t)}$$

Model uncertainty

So far, our algorithms have considered *fixed* parameter spaces. What if, however, there is some uncertainty as to which model is the “correct” one from a set $\{\mathcal{M}_1, \dots, \mathcal{M}_k\}$? We can include the model as an additional parameter to be estimated

$$\pi(\boldsymbol{\theta}_{\mathcal{M}}, \mathcal{M} \mid \mathbf{x}) \propto L(\mathbf{x} \mid \boldsymbol{\theta}_{\mathcal{M}}, \mathcal{M})p(\boldsymbol{\theta}_{\mathcal{M}}, \mathcal{M})p(\mathcal{M})$$

The idea is then to construct a Markov chain that is able to move between different models, and has stationary distribution π , to generate samples $\{\mathcal{M}^{(t)}, \boldsymbol{\theta}_{\mathcal{M}^{(t)}}^{(t)}\}_{t=1}^T$.

Using these samples, here are two examples of statistics we might be interested in calculating:

1. The posterior probability of each model

$$\pi(\mathcal{M}_i \mid \mathbf{x}) = \frac{L(\mathbf{x} \mid \mathcal{M}_i)p(\mathcal{M}_i)}{\sum_{i=1}^k L(\mathbf{x} \mid \mathcal{M}_i)p(\mathcal{M}_i)}$$

where

$$L(\mathbf{x} \mid \mathcal{M}_i) = \int L(\mathbf{x} \mid \mathcal{M}_i, \boldsymbol{\theta}_{\mathcal{M}_i})p(\boldsymbol{\theta}_{\mathcal{M}_i} \mid \mathcal{M}_i) d\boldsymbol{\theta}_{\mathcal{M}_i}$$

This can be calculated, using $\{\mathcal{M}^{(t)}, \boldsymbol{\theta}_{\mathcal{M}^{(t)}}\}_{t=1}^T$, as

$$\hat{\pi}(\mathcal{M}_i | \mathbf{x}) = \frac{1}{T} \sum_{t=1}^T \mathbb{I}_{\{\mathcal{M}^{(t)} = \mathcal{M}_i\}}$$

2. Estimates concerning the parameters themselves, based on the distribution

$$\pi(\boldsymbol{\theta} | \mathbf{x}) = \sum_{i=1}^k \pi(\boldsymbol{\theta} | \mathbf{x}, \mathcal{M}_i) \pi(\mathcal{M}_i | \mathbf{x})$$

We can calculate such statistics by simply ignoring the model part of our Markov chain. For example

$$\begin{aligned} \hat{\boldsymbol{\theta}} &= \mathbb{E}(\boldsymbol{\theta}) = \frac{1}{T} \sum_{i=1}^T \boldsymbol{\theta}^{(i)} \\ \text{var}(\boldsymbol{\theta}) &= \frac{1}{T} \sum_{i=1}^T (\boldsymbol{\theta}^{(i)} - \hat{\boldsymbol{\theta}})^2 \end{aligned}$$

Of course, this only works if $\boldsymbol{\theta}$ is a statistic that is common to all the models (for example, the mean – see the example below).

The algorithm we use to generate our desired Markov Chain is called the **reversible jump algorithm** (some written as RJ-MCMC).

Reversible-Jump Monte Carlo:

1. Given a model \mathcal{M} , update the parameters $\boldsymbol{\theta}_{\mathcal{M}}$ using MH or GS.
2. Reversible jump (RJ)-step: with probability $P(\mathcal{M} \rightarrow \mathcal{M}')$ [which must be chosen], propose to replace the model \mathcal{M} with a new model \mathcal{M}' . This proposal can be accepted or rejected.

The RJ step is complicated because it involves proposing new parameters $\boldsymbol{\theta}_{\mathcal{M}'}$ for the new model. Imagine the move we are proposing is

$$\{\mathcal{M}, \boldsymbol{\theta}_{\mathcal{M}}\} \rightarrow \{\mathcal{M}', \boldsymbol{\theta}_{\mathcal{M}'}\}$$

Our approach will be to choose

$$(\boldsymbol{\theta}_{\mathcal{M}'}, \mathbf{u}') = g_{\mathcal{M}, \mathcal{M}'}(\boldsymbol{\theta}_{\mathcal{M}}, \mathbf{u})$$

Where

- g is a known, deterministic, bijective function⁵ between $(\boldsymbol{\theta}_{\mathcal{M}'}, \mathbf{u})$ and $(\boldsymbol{\theta}_{\mathcal{M}}, \mathbf{u}')$.
- \mathbf{u} and \mathbf{u}' are vectors chosen such that

$$\dim \mathbf{u} + \dim(\boldsymbol{\theta}_{\mathcal{M}}) = \dim \mathbf{u}' + \dim(\boldsymbol{\theta}_{\mathcal{M}'})$$

In other words, it “matches” the dimensions of the parameters. Either one or both of \mathbf{u} and \mathbf{u}' will be 0, depending on the size of the models.

If $\dim(\boldsymbol{\theta}_{\mathcal{M}'}) > \dim(\boldsymbol{\theta}_{\mathcal{M}})$, \mathbf{u} is nonzero. We choose it by sampling from a proposal density $q_{\mathcal{M}, \mathcal{M}'}(\mathbf{u})$

We then accept the move with probability $\alpha(\boldsymbol{\theta}_{\mathcal{M}}, \boldsymbol{\theta}_{\mathcal{M}'}) = \min\{1, A\}$

$$A = \underbrace{\frac{\pi(\mathcal{M}', \boldsymbol{\theta}_{\mathcal{M}'} | \mathbf{x})}{\pi(\mathcal{M}, \boldsymbol{\theta}_{\mathcal{M}} | \mathbf{x})}}_{\substack{\text{Model ratio} \\ (= \frac{L(\mathbf{x} | \mathcal{M}', \boldsymbol{\theta}_{\mathcal{M}'}) p(\boldsymbol{\theta}_{\mathcal{M}' | \mathcal{M}')} p(\mathcal{M}')}{L(\mathbf{x} | \mathcal{M}, \boldsymbol{\theta}_{\mathcal{M}}) p(\boldsymbol{\theta}_{\mathcal{M} | \mathcal{M}}) p(\mathcal{M})})}} \cdot \underbrace{\frac{P(\mathcal{M}' \rightarrow \mathcal{M}) q(\mathbf{u}')}{P(\mathcal{M} \rightarrow \mathcal{M}') q(\mathbf{u})} \left| \frac{\partial g(\boldsymbol{\theta}_{\mathcal{M}}, \mathbf{u})}{\partial(\boldsymbol{\theta}_{\mathcal{M}}, \mathbf{u})} \right|}_{\text{Proposal ratio}}$$

Note that

$$\left(\frac{\partial g(\boldsymbol{\theta}_{\mathcal{M}}, \mathbf{u})}{\partial(\boldsymbol{\theta}_{\mathcal{M}}, \mathbf{u})} \right)_{ij} = \left(\frac{\partial(\boldsymbol{\theta}_{\mathcal{M}'}, \mathbf{u}')}{\partial(\boldsymbol{\theta}_{\mathcal{M}}, \mathbf{u})} \right)_{ij} = \frac{\partial((\boldsymbol{\theta}_{\mathcal{M}'}, \mathbf{u}')_i)}{\partial((\boldsymbol{\theta}_{\mathcal{M}}, \mathbf{u})_j)}$$

Where $(\boldsymbol{\theta}_{\mathcal{M}}, \mathbf{u})$ is the vector formed by combining $\boldsymbol{\theta}_{\mathcal{M}}$ and \mathbf{u} .

Further note that the reverse move from $\{\mathcal{M}', \boldsymbol{\theta}_{\mathcal{M}'}\} \rightarrow \{\mathcal{M}, \boldsymbol{\theta}_{\mathcal{M}}\}$ is fully defined by g^{-1} and $\alpha(\boldsymbol{\theta}_{\mathcal{M}'}, \boldsymbol{\theta}_{\mathcal{M}}) = \min\{1, A^{-1}\}$.

EXAMPLE: Suppose that we observe data \mathbf{x} which are IID, but where the distribution is unknown; either $\text{Exp}(\lambda)$ or $\text{Gamma}(\alpha, \beta)$, with all parameters unknown⁶. We let $\mathbb{P}(\text{Exp})$ and $\mathbb{P}(\text{G})$ be the prior probability on each model, and let the priors on the parameters be

$$\lambda \sim \text{Gamma}(a_1, b_1) \quad \alpha \sim \text{Gamma}(a_2, b_2) \quad \beta \sim \text{Gamma}(a_3, b_3)$$

Let’s find posterior distributions

Exponential model

⁵ A bijective function f from a set X to a set Y has the property that for every y in Y , there is exactly one x in X such that $f(x) = y$ and no unmapped element exists in either X or Y .

⁶ This example is really quite silly, because $\text{Gamma}(1, \lambda) \sim \text{Exp}(\lambda)$, but it’s useful for demonstration.

[Reminder: we need all the constants in the original distribution, because λ is something we'll want the distribution of!]

$$L(\mathbf{x} \mid \lambda, \text{Exp}) = \prod_{i=1}^n \lambda \exp(-\lambda x_i) = \lambda^n \exp(-\lambda x_+)$$

$$\Downarrow$$

$$\pi(\lambda \mid \mathbf{x}) \propto L(\mathbf{x} \mid \lambda, \text{Exp}) p(\lambda) \propto \lambda^{n+a_1-1} e^{-\lambda x_+ - \lambda b_1}$$

$$\boxed{\lambda \mid \mathbf{x} \sim \text{Gamma}(n + a_1, x_+ + b_1)}$$

Gamma Model

[Reminder: we need all the constants in the original distribution, because α and β are things we'll want the distribution of later! The Bayes' factors, however, integrate out the relevant variables by definition, and so can be ignored]

$$L(\mathbf{x} \mid \alpha, \beta, \text{G}) = \prod_{i=1}^n \frac{\beta^\alpha}{\Gamma(\alpha)} x_i^{\alpha-1} \exp(-\beta x_i) = \left(\frac{\beta^\alpha}{\Gamma(\alpha)}\right)^n x_i^{n(\alpha-1)} \exp(-\beta x_+)$$

$$\Downarrow$$

$$\pi(\beta \mid \mathbf{x}, \alpha) \propto L(\mathbf{x} \mid \alpha, \beta, \text{G}) p(\beta) \propto \beta^{\alpha n + a_3 - 1} \exp(-\beta x_+ - b_3 \beta)$$

$$\boxed{\beta \mid \mathbf{x}, \alpha \sim \text{Gamma}(\alpha n + a_3, x_+ + b_3)}$$

$$\Downarrow$$

$$\pi(\alpha \mid \mathbf{x}, \beta) \propto L(\mathbf{x} \mid \alpha, \beta, \text{G}) p(\alpha)$$

$$\boxed{\pi(\alpha \mid \mathbf{x}, \beta) \propto \left(\frac{\beta^\alpha}{\Gamma(\alpha)}\right)^n x_i^{n(\alpha-1)} \alpha^{a_2-1} \exp(-b_2 \alpha)}$$

At each step, the first two distributions can be sampled from using a GS or the MH algorithm. We given an example for α and β . The latter can be sampled from using a simple Gibbs sampler. The former needs MH:

- Start with a value of $\alpha^{(0)}$
- Sample $\beta^{(t+1)}$ from $\text{Gamma}(\alpha^{(t)}n + a_3, x_+ + b_3)$
- Simulate ϕ from $U[\frac{3}{4}\alpha, \frac{4}{3}\alpha]$ (this is our proposal density).
- Calculate the acceptance probability

$$A(\alpha^{(t)}, \phi) = \min\left(1, \frac{\pi(\phi \mid \beta^{(t)}, \mathbf{x}) \frac{1}{(\frac{4}{3}-\frac{3}{4})\phi} \mathbb{I}_{\alpha^{(t)} \in [\frac{3}{4}\phi, \frac{4}{3}\phi]}}{\pi(\alpha^{(t)} \mid \beta^{(t)}, \mathbf{x}) \frac{1}{(\frac{4}{3}-\frac{3}{4})\alpha^{(t)}} \mathbb{I}_{\phi \in [\frac{3}{4}\alpha^{(t)}, \frac{4}{3}\alpha^{(t)}]}}\right)$$

With probability A , set $\alpha^{(t+1)} = \phi$, and with probability $1 - A$, set $\alpha^{(t+1)} = \alpha^{(t)}$.

We can now consider the RJ step. In each case, we will “propose” to move with probability 1 (so $P(\text{Exp} \rightarrow \text{Gamma}) = P(\text{Gamma} \rightarrow \text{Exp}) = 1$).

Let us now consider the move from Exp to Gamma. We must deal with the augmentation of the parameter space, from $(\lambda) \rightarrow (\alpha, \beta)$.

- The target parameter space is larger, so we set $u' = 0$, and we sample $u \sim \text{Gamma}(\gamma, \delta)$.
- We define our function g so that the mean of the original Exp distribution and of the resulting Gamma distribution is the same. This gives

$$g(\lambda, u) \rightarrow (\alpha = u, \beta = \lambda u)$$

- We then have

$$\left| \frac{\partial(\alpha, \beta)}{\partial(\lambda, u)} \right| = \begin{vmatrix} \partial\alpha / \partial\lambda & \partial\alpha / \partial u \\ \partial\beta / \partial\lambda & \partial\beta / \partial u \end{vmatrix} = \begin{vmatrix} 0 & 1 \\ u & \lambda \end{vmatrix} = u$$

The probability of accepting the model move is then given by $\min\{1, A\}$, where

$$\begin{aligned} A &= \frac{\pi(\text{G}, \alpha, \beta \mid \mathbf{x})}{\pi(\text{Exp}, \lambda \mid \mathbf{x})} \frac{P(\text{G} \rightarrow \text{Exp})}{P(\text{Exp} \rightarrow \text{G})} \frac{\left| \frac{\partial(\alpha, \beta)}{\partial(\lambda, u)} \right|}{q(u)} \\ &= \frac{L(\mathbf{x} \mid \text{G}, \alpha, \beta) p(\alpha) p(\beta) \mathbb{P}(\text{G})}{L(\mathbf{x} \mid \lambda, \text{Exp}) p(\lambda) \mathbb{P}(\text{Exp})} \frac{u}{\frac{1}{\Gamma(\gamma)} \delta^\gamma u^{\gamma-1} e^{-\delta u}} \end{aligned}$$

The reverse step is defined by $(\lambda, u) = g^{-1}(\alpha, \beta) \Rightarrow u = \alpha, \lambda = \beta / \alpha$ and probability A^{-1} .

EXAMPLE: Suppose that we observe data \mathbf{x} which are IID, but where the distribution is unknown; either $\text{Exp}(\lambda)$ or $\text{Pareto}(\alpha, \beta)$, with all parameters unknown. We let $\mathbb{P}(\text{Exp})$ and $\mathbb{P}(\text{Par})$ be the prior probability on each model, and let the priors on the parameters be

$$\lambda \sim \text{Gamma}(a, b) \quad \alpha \sim \text{Exp}(\mu) \quad \beta \sim \text{Exp}(\nu)$$

Let's find posterior distributions

Exponential model

$$L(\mathbf{x} \mid \lambda, \text{Exp}) = \prod_{i=1}^n \lambda \exp(-\lambda x_i) = \lambda^n \exp(-\lambda x_+)$$

$$\Downarrow$$

$$\pi(\lambda \mid \mathbf{x}) \propto L(\mathbf{x} \mid \lambda, \text{Exp})p(\lambda) \propto \lambda^{n+a-1} e^{-\lambda x_+ - \lambda b}$$

$$\boxed{\lambda \mid \mathbf{x} \sim \text{Gamma}(n + a, x_+ + b)}$$

Pareto Model

The Pareto distribution comes with an additional complication, in that it is imperative that $0 < \alpha \leq x$. Thus

$$L(\mathbf{x} \mid \alpha, \beta, \text{Par}) = \prod_{i=1}^n \beta \alpha^\beta x_i^{-(\beta+1)} = \beta^n \alpha^{n\beta} \prod_{i=1}^n x_i^{-(\beta+1)}$$

$$\Downarrow$$

$$\pi(\alpha \mid \beta, \mathbf{x}) \propto \beta^n \alpha^{n\beta} \left(\prod_{i=1}^n x_i^{-(\beta+1)} \right) \mu e^{-\mu \alpha} \mathbb{I}_{\{\alpha \in (0, x_{\min}]\}}$$

$$\boxed{\pi(\alpha \mid \beta, \mathbf{x}) \propto \alpha^{n\beta} e^{-\mu \alpha} \mathbb{I}_{\{\alpha \in (0, x_{\min}]\}}}$$

$$\Downarrow$$

$$\pi(\beta \mid \alpha, \mathbf{x}) \propto \beta^n \alpha^{n\beta} \left(\prod_{i=1}^n x_i^{-(\beta+1)} \right) \nu e^{-\nu \beta}$$

$$\pi(\beta \mid \alpha, \mathbf{x}) \propto \beta^n \exp(-\nu \beta + n\beta \log \alpha - \beta \sum_{i=1}^n \log x_i)$$

$$\boxed{\beta \mid \alpha, \mathbf{x} \sim \text{Gamma}(n + 1, \nu - n \log \alpha + \sum_{i=1}^n \log x_i)}$$

Once again, the first distribution is not in standard form, and therefore requires an MH sampler. We can use a symmetric random walk MH algorithm as follows:

- Given an α , generate $\alpha' = \alpha + u$, where $u \sim U[0, \sigma^2]$
- Automatically reject α' if $\alpha' \notin (0, x_{\min}]$
- Otherwise, accept it with probability

$$\min \left(1, \frac{\pi(\alpha' \mid \beta, \mathbf{x})}{\pi(\alpha \mid \beta, \mathbf{x})} \right) = \min \left(1, \frac{\alpha'^{n\beta} \exp(-\mu \alpha')}{\alpha^{n\beta} \exp(-\mu \alpha)} \right)$$

We can now consider the RJ step. In each case, we will “propose” to move with probability 1 (so $P(\text{Exp} \rightarrow \text{Par}) = P(\text{Par} \rightarrow \text{Exp}) = 1$).

Let us now consider the move from Exp to Pareto. We must deal with the augmentation of the parameter space, from $(\lambda) \rightarrow (\alpha, \beta)$.

- The target parameter space is larger, so we set $u' = 0$, and we sample $u \sim \text{Exp}(\gamma)$.
- We define our function g so that the mean of the original Exp distribution and of the resulting Pareto distribution is the same. This gives

$$g(\lambda, u) \rightarrow \left(\alpha = \frac{u}{\lambda(u+1)}, \beta = u + 1 \right)$$

- We then have

$$\begin{aligned} \left| \frac{\partial(\alpha, \beta)}{\partial(\lambda, u)} \right| &= \begin{vmatrix} \partial\alpha / \partial\lambda & \partial\alpha / \partial u \\ \partial\beta / \partial\lambda & \partial\beta / \partial u \end{vmatrix} \\ &= \begin{vmatrix} -\frac{u}{\lambda^2(u+1)} & \frac{1}{\lambda(u+1)} - \frac{u}{\lambda(u+1)^2} \\ 0 & 1 \end{vmatrix} \\ &= \frac{u}{\lambda^2(u+1)} \end{aligned}$$

The probability of accepting the model change is then given by $\min\{1, A\}$, where

$$\begin{aligned} A &= \frac{\pi(\text{Par}, \alpha, \beta | \mathbf{x})}{\pi(\text{Exp}, \lambda | \mathbf{x})} \frac{P(\text{Par} \rightarrow \text{Exp})}{P(\text{Exp} \rightarrow \text{Par})} \frac{\left| \frac{\partial(\alpha, \beta)}{\partial(\lambda, u)} \right|}{q(u)} \\ &= \frac{\pi(\alpha, \beta | \mathbf{x}, \text{Par}) p(\text{Par})}{\pi(\lambda | \mathbf{x}, \text{Exp}) p(\text{Exp})} \frac{1}{\phi\left(u/\sqrt{\sigma^2}\right)} \frac{u}{\lambda^2(u+1)} \\ &= \frac{L(\mathbf{x} | \alpha, \beta, \text{Par}) p(\alpha) p(\beta) p(\text{Par})}{L(\mathbf{x} | \lambda, \text{Exp}) p(\lambda) p(\text{Exp})} \frac{1}{\phi\left(u/\sqrt{\sigma^2}\right)} \frac{u}{\lambda^2(u+1)} \end{aligned}$$

Sequential Importance Sampling

Sequential importance sampling is an alternative to MCMC for sampling from high-dimensional distributions. It helps with the problem of estimating good proposal/importance distributions by building them sequentially.

To motivate the concept, first consider a target distribution $\pi(\boldsymbol{\theta})$ and a proposal distribution $g(\boldsymbol{\theta})$. Denoting $\boldsymbol{\theta}_{(-j)} = (\theta_1, \dots, \theta_j, \dots, \theta_p)$ We can write these as

$$g(\boldsymbol{\theta}) = g_1(\theta_1)g_2(\theta_2 | \theta_1) \cdots g_p(\theta_p | \boldsymbol{\theta}_{(-p)})$$

$$\pi(\boldsymbol{\theta}) = \pi(\theta_1)\pi(\theta_2 | \theta_1) \cdots \pi(\theta_p | \boldsymbol{\theta}_{(-p)})$$

And we can write the importance weights as

$$w(\boldsymbol{\theta}) = \frac{\pi(\theta_1)\pi(\theta_2 | \theta_1) \cdots \pi(\theta_p | \boldsymbol{\theta}_{(-p)})}{g_1(\theta_1)g_2(\theta_2 | \theta_1) \cdots g_p(\theta_p | \boldsymbol{\theta}_{(-p)})}$$

Writing $\boldsymbol{\theta}_k = (\theta_1, \dots, \theta_k)$, we can define a *partial weight*

$$w_k(\boldsymbol{\theta}_k) = w_{k-1}(\boldsymbol{\theta}_{k-1}) \frac{\pi(\theta_k | \boldsymbol{\theta}_{k-1})}{g_k(\theta_k | \boldsymbol{\theta}_{k-1})}$$

When we include the entire vector, $w_p(\boldsymbol{\theta}_p) = w(\boldsymbol{\theta})$. Using this method is advantageous for two reasons

- We can stop generating further components if the partial weight gets too small (this will, however, lead to bias).
- We can use the marginal distribution $\pi(\theta_k | \boldsymbol{\theta}_{k-1})$ to help us in designing $g_k(\theta_k | \boldsymbol{\theta}_{k-1})$.

The only problem is that the decomposition of π is very difficult; it requires

$$\pi(\boldsymbol{\theta}_k) = \int \pi(\boldsymbol{\theta}) d\theta_{k+1} \cdots d\theta_p$$

Which is as difficult (or harder than) the initial problem.

Sequential importance sampling: Suppose we can find a sequence of “auxiliary distributions”, which need not be normalised, such that

$$\pi_k(\boldsymbol{\theta}_k) \approx \pi(\boldsymbol{\theta})_k \quad k = 1, \dots, p - 1$$

$$\pi_p(\boldsymbol{\theta}_p) = \pi(\boldsymbol{\theta})_p$$

The SIS method is then defined as the following recursive procedure

1. Draw $\theta_k \sim g_k(\theta_k | \boldsymbol{\theta}_{k-1})$, and let $\boldsymbol{\theta}_k = (\boldsymbol{\theta}_{k-1}, \theta_k)$.
2. Compute the *incremental weight*

$$u_k = \frac{\pi_k(\boldsymbol{\theta}_k)}{\pi_{k-1}(\boldsymbol{\theta}_{k-1}) g_k(\theta_k | \boldsymbol{\theta}_{k-1})} \cdot 1$$

3. Let $w_k = w_{k-1} u_k$

Note: An important special case arises when we can build g_k using π_k ; ie:

$$g_k(\boldsymbol{\theta}_k | \boldsymbol{\theta}_{k-1}) = \pi_k(\boldsymbol{\theta}_k | \boldsymbol{\theta}_{k-1})$$

in that case,

$$u_k = \frac{\pi_k(\boldsymbol{\theta}_{k-1})}{\pi_{k-1}(\boldsymbol{\theta}_{k-1})}$$

We often obtain several samples $\boldsymbol{\theta}^{(1)}, \dots, \boldsymbol{\theta}^{(T)}$ in parallel. Each of the partial samples $\mathcal{S}_k = \{\boldsymbol{\theta}_k^{(t)}\}_{t=1}^T$ are called *particles* or *streams* when properly weighed by the collection of weights $\mathcal{W}_k = \{w_k^{(t)}\}_{t=1}^T$.

As k increases, the variance of the importance weights also increases, which decreases the overall effective sample size. To fix this problem, we can use *resampling*

Resampling: Periodically, randomly, or dynamically (for example, when the ESS) is low, perform the following two steps:

1. Sample a new set of steams \mathcal{S}'_k from \mathcal{S}_k , with replacement and with weights \mathcal{W}_k .
2. Assign equal weights $\sum_{t=1}^T w_k^{(t)} / T$ to each of the steams in \mathcal{S}'_k . (It is also common to set the weights to $1/T$).

Notes: Resampling can result in few unique; this is known as *particle depletion*. *Enrichment* or *diversification* methods involving kernel smoothing and MCMC have been developed as a remedy.

SIS is part of a broader class of Sequential Monte Carlo (SMC) methods, which are popular for inference in state space models, and are commonly encountered when dealing with time series data.

EXAMPLE: Suppose the auxiliary distributions are only known up to a normalising constant. Denote the normalised distributions by π_k and the un-normalised distributions by

$$q_k(\boldsymbol{\theta}_k) = Z_k \pi_k(\boldsymbol{\theta}_k) \quad k = 1, \dots, p$$

The incremental weights take the form

$$u_k = \frac{q_k(\boldsymbol{\theta}_k)}{q_{k-1}(\boldsymbol{\theta}_{k-1})g_k(\boldsymbol{\theta}_k | \boldsymbol{\theta}_{k-1})} = \frac{Z_k \pi_k(\boldsymbol{\theta}_k)}{Z_{k-1} \pi_{k-1}(\boldsymbol{\theta}_{k-1})g_k(\boldsymbol{\theta}_k | \boldsymbol{\theta}_{k-1})}$$

and final weight takes the form

$$w_p = \prod_{i=2}^p u_i = \frac{Z_p}{Z_1} \frac{\pi_p(\boldsymbol{\theta}_p)}{g_1(\boldsymbol{\theta}_1) \cdots g_p(\boldsymbol{\theta}_p | \boldsymbol{\theta}_{p-1})}$$

Thus, the sample average of $\{w_p^{(t)}\}_{t=1}^T$ gives an unbiased estimate of $\mathbb{E}(w_p) = Z_p / Z_1$.

Classical Inference

Classical inference involves maximising a likelihood $L(\mathbf{x} | \boldsymbol{\theta})$ to obtain estimates $\hat{\boldsymbol{\theta}}$ for the parameters. It will often be more convenient to consider *minimising* $-L(\mathbf{x} | \boldsymbol{\theta})$.

Simulated Annealing

Sometimes, this can be done analytically, but often, numerical algorithms are required. Some, called *absolute descent* algorithms try to move downhill until they can get no lower. Unfortunately, these algorithms are susceptible to getting stuck at subsidiary maxima. Consider, instead, an algorithm that proposes a move from $\boldsymbol{\theta} \rightarrow \boldsymbol{\phi}$ such that

- If $f(\boldsymbol{\phi}) \leq f(\boldsymbol{\theta})$, the move is accepted (absolute descent)
- If $f(\boldsymbol{\phi}) > f(\boldsymbol{\theta})$, the move is accepted with probability α

The problem with this method is that it does not present any obvious place to “stop”. One approach would be to run the algorithm for a fixed time (but then it’s unclear when to stop).

Another approach, which is adopted by *simulated annealing*, is to decrease α at each step. Eventually, the system freezes, to what we would hope is the global optimum.

More specifically, suppose we want to minimise the function $f(\boldsymbol{\theta})$, and let

$$b_T(\boldsymbol{\theta}) = \frac{\exp\{-f(\boldsymbol{\theta})/T\}}{\int \exp\{-f(\boldsymbol{\theta})/T\} d\boldsymbol{\theta}} \propto \exp\left\{-\frac{f(\boldsymbol{\theta})}{T}\right\}$$

The parameter T is called the *temperature*. Clearly, b_T is a distribution that favours smaller values of $f(\boldsymbol{\theta})$. As $T \rightarrow 0$, smaller values of $f(\boldsymbol{\theta})$ are increasingly preferred.

Ideally, we would like to sample from b_0 , but this is a point mass of the unknown minimising value $\hat{\boldsymbol{\theta}}$. Instead, simulated annealing works by simulating from a series of b_T for a decreasing sequence of temperatures:

Simulated Annealing (SA):

1. Take an initial temperature T_0 and a starting value $\boldsymbol{\theta}_0$.

2. Propose a new state ϕ with density $q(\theta, \phi)$ [a typical choice is random-walk Metropolis, with $\phi = \theta + z$].
3. Accept the move with probability $\alpha(\theta, \phi) = \min\{1, A\}$, with

$$A = \frac{b_T(\phi)q(\phi, \theta)}{b_T(\theta)q(\theta, \phi)}$$
4. Repeat steps 2 and 3 until the chain reaches equilibrium.
5. Check for a stopping criterion. A common criterion is to stop if no moves were accepted at this temperature.
6. Lower the temperature T and return to step 2.

Notes: When maximising a likelihood, we typically set $f(\theta) = -L(\mathbf{x} | \theta)$ or $f(\theta) = -\log L(\mathbf{x} | \theta)$. In the latter case, $b_T(\theta) \propto [L(\mathbf{x} | \theta)]^{1/T}$.

The overall Markov chain is *inhomogeneous* because the target distribution changes over time.

The ideal proposal distribution q for the MH sampler still depend on the temperature T – the lower T , the closer we'd expect to be to the minimum, and the lower the variance of q . For this reason, Gibbs Sampling is preferred to standard distributions, since no proposal distribution is needed.

EXAMPLE: Consider trying to find the mode of a distribution comprised of a mixture of normals:

$$\begin{aligned} f(m) &= -\log\left(L(m | \mu_1, \sigma_1^2, \mu_2, \sigma_2^2)\right) \\ &= -\log\left(0.6N(m | \mu_1 = -8, \sigma_1^2 = 0.5^2) \right. \\ &\quad \left. + 0.4N(m | \mu_2 = 8, \sigma_2^2 = 0.9^2)\right) \end{aligned}$$

And so

$$b_T(m) = \left[0.6\phi\left(\frac{m - \mu_1}{\sigma_1}\right) + 0.4\phi\left(\frac{m - \mu_2}{\sigma_2}\right) \right]^{1/T}$$

For fixed T , sampling from b_T will require MH.

EXAMPLE: Consider, now, the AR(k) example. If, once again, we let $f(\boldsymbol{\theta}) = -\log L(\mathbf{x} | \boldsymbol{\theta})$, we get:

$$\begin{aligned} b_T(\mathbf{a}, \sigma^2) &\propto (\sigma^{-2})^{\frac{N-k}{2T}} \exp\left(-\frac{1}{2T\sigma^2} (\mathbf{x} - X_k \mathbf{a})^2\right) \\ &= (\sigma^{-2})^{\frac{N-k}{2T}} \exp\left(-\frac{1}{2T\sigma^2} \sum_{t=k+1}^N (x_t - \sum_{r=1}^k a_r x_{t-r})^2\right) \end{aligned}$$

Finding the full conditionals

- **For σ^2**

$$b_T(\mathbf{a}, \sigma^2) \propto (\sigma^{-2})^{\frac{N-k}{2T}} \exp\left(-\frac{1}{\sigma^2} \frac{1}{2T} \sum_{t=k+1}^N (x_t - \sum_{r=1}^k a_r x_{t-r})^2\right)$$

And so

$$b_T(\mathbf{a}, \sigma^2) \sim \text{InvGamma}\left(\frac{N-k}{2T} + 1, \frac{1}{2T} \sum_{t=k+1}^N (x_t - \sum_{r=1}^k a_r x_{t-r})^2\right)$$

- **And for \mathbf{a}**

$$\begin{aligned} b_T(\mathbf{a} | \sigma^2) &\propto \exp\left(-\frac{1}{2T\sigma^2} (\mathbf{x} - X_k \mathbf{a})^T (\mathbf{x} - X_k \mathbf{a})\right) \\ &\propto \exp\left(-\frac{1}{2} \left\{ \mathbf{a}^T \frac{X_k^T X_k}{T\sigma^2} \mathbf{a} - 2 \frac{\mathbf{a}^T X_k^T \mathbf{x}}{T\sigma^2} \right\}\right) \\ &\propto \exp\left(-\frac{1}{2} \left\{ \mathbf{a}^T \Sigma^{-1} \mathbf{a} - 2 \mathbf{a}^T \Sigma^{-1} \boldsymbol{\mu} \right\}\right) \\ &\propto \exp\left(-\frac{1}{2} (\mathbf{a} - \boldsymbol{\mu})^T \Sigma^{-1} (\mathbf{a} - \boldsymbol{\mu})\right) \\ &\sim N\left(\boldsymbol{\mu}, \frac{\Sigma^{-1}}{T\sigma^2}\right) \end{aligned}$$

Where

$$\Sigma^{-1} = \frac{X_k^T X_k}{T\sigma^2} \qquad \boldsymbol{\mu} = \frac{\Sigma X_k^T \mathbf{x}}{T\sigma^2}$$

Note that as $T \rightarrow 0$, the variance tends to 0, but the mean is unchanged.

Both can be sampled using a Gibbs sampler.

Model Selection Using Simulated Annealing

We can also use simulated annealing to do model selection. For example, to use Akaike's information criterion, we set

$$f(\boldsymbol{\theta}_{\mathcal{M}}, \mathcal{M}) = -2 \log L_{\mathcal{M}}(\mathbf{x} \mid \boldsymbol{\theta}_{\mathcal{M}}) + 2 \dim(\mathcal{M})$$

$$g_T(\boldsymbol{\theta}_{\mathcal{M}}, \mathcal{M}) = \exp(-f(\boldsymbol{\theta}_{\mathcal{M}}, \mathcal{M}) / T)$$

Of course, we now require inter-dimensional jumps, and so RJ-MCMC must be used.

EXAMPLE: Consider, now, the AR(k) example. Use

$$f(\mathbf{a}_k, \sigma^2, k) = -\log L_k(\mathbf{x} \mid \mathbf{a}_k, \sigma^2) + k$$

$$b_T(\mathbf{a}, \sigma^2, k) \propto \left(\sigma^{-2}\right)^{\frac{N-k}{T}} \exp\left\{-\frac{1}{2\sigma^2 T}(\mathbf{x} - X_k \mathbf{a})^2 - \frac{k}{T}\right\}$$

Now suppose we wish to jump from $k \rightarrow k+1$. In terms of parameters, we're going from

$$g\left\{(\mathbf{a}_k, u), \sigma^2\right\} = \left\{\mathbf{a}_{k+1}, \sigma^2\right\}$$

And we generate $u \sim q$. The Jacobian is the identity, since the mapping g is the identity.

The proposed move is then accepted with probability $\alpha = \min\{1, A\}$, where

$$A = \frac{b_T(\mathbf{a}_{k+1}, \sigma^2, k+1)}{b_T(\mathbf{a}_k, \sigma^2, k)} \frac{P(k+1 \rightarrow k)}{P(k \rightarrow k+1)q(u)}$$

Now, we saw in the previous example how to work out $b_T(\mathbf{a}_{k+1} \mid \mathbf{a}_k, \sigma^2, k+1)$; it's a normal distribution, given by

$$b_T(\mathbf{a}_{k+1} \mid \mathbf{a}_k, \sigma^2, k+1) = \frac{b_T(\mathbf{a}_{k+1}, \sigma^2, k+1)}{\int b_T(\mathbf{a}_{k+1}, \sigma^2, k+1) da_{k+1}}$$

It seems sensible to use this as our proposal distribution q . Using that and $P(k+1 \rightarrow k) = P(k \rightarrow k+1)$, the acceptance probability reduces to

$$A = \frac{\int b_T(\mathbf{a}_{k+1}, \sigma^2, k+1) da_{k+1}}{b_T(\mathbf{a}_k, \sigma^2, k)}$$

ie: the ratio of the marginal distributions of the unchanged parameters under each model.

Expectation Maximisation

Expectation minimisation is a method that is used when the likelihood is not defined explicitly, but is known in the form

$$f(\boldsymbol{\theta}) = \int g(\boldsymbol{\theta}, \mathbf{z}) \, d\mathbf{z} \quad \text{or} \quad f(\boldsymbol{\theta}) = \int g_1(\boldsymbol{\theta} | \mathbf{z})g_2(\mathbf{z}) \, d\mathbf{z}$$

An obvious situation in which this occurs is where we have known data, contained in the vector \mathbf{x} , and missing data, contained in the vector \mathbf{z} . In that case, the likelihood we want to maximise is

$$L(\mathbf{x} | \boldsymbol{\theta}) = \int L(\mathbf{x} | \mathbf{z}, \boldsymbol{\theta})f(\mathbf{z}) \, d\mathbf{z}$$

Expectation Maximisation (EM):

1. *E-Step* – given \mathbf{x} , calculate the expectation of the complete data log-likelihood as a function of the current estimate of $\boldsymbol{\theta}$

$$\begin{aligned} Q(\boldsymbol{\theta} | \boldsymbol{\theta}^{(t)}) &= \mathbb{E} \left\{ \ell(\mathbf{x}, \mathbf{z} | \boldsymbol{\theta}) | \mathbf{x}, \boldsymbol{\theta}^{(t)} \right\} \\ &= \mathbb{E}_{\mathbf{z} | \mathbf{x}, \boldsymbol{\theta}^{(t)}} \left\{ \ell(\mathbf{x}, \mathbf{z} | \boldsymbol{\theta}) \right\} \\ &= \int \ell(\mathbf{x}, \mathbf{z} | \boldsymbol{\theta})f(\mathbf{z} | \mathbf{x}, \boldsymbol{\theta} = \boldsymbol{\theta}^{(t)}) \, d\mathbf{z} \end{aligned}$$

if the log-likelihood is linear in the joint sufficient statistics of $[\mathbf{x}, \mathbf{z}]$, then this step simply involves finding the expectation of \mathbf{z} given \mathbf{x} and $\boldsymbol{\theta}^{(t)}$ and feeding it into ℓ .

2. *M-Step* – find $\boldsymbol{\theta}^{(t+1)}$ with maximises $Q(\boldsymbol{\theta} | \boldsymbol{\theta}^{(t)})$

Theorem: Every step of the EM algorithm increases the log likelihood. That is

$$\ell(\mathbf{x} | \boldsymbol{\theta}^{(t+1)}) \geq \ell(\mathbf{x} | \boldsymbol{\theta}^{(t)})$$

with equality if and only if

$$Q(\boldsymbol{\theta}^{(t+1)} | \boldsymbol{\theta}^{(t)}) = Q(\boldsymbol{\theta}^{(t)} | \boldsymbol{\theta}^{(t)})$$

Proof: The likelihood of the complete data can be factorised as

$$L(\mathbf{x}, \mathbf{z} | \boldsymbol{\theta}) = L(\mathbf{x} | \boldsymbol{\theta})f(\mathbf{z} | \mathbf{x}, \boldsymbol{\theta})$$

(To understand why this expression makes sense, we can view L as a sort of “prior on \mathbf{x} ”). And so

$$\begin{aligned}\ell(\mathbf{x}, \mathbf{z} \mid \boldsymbol{\theta}) &= \ell(\mathbf{x} \mid \boldsymbol{\theta}) + \log f(\mathbf{z} \mid \mathbf{x}, \boldsymbol{\theta}) \\ \ell(\mathbf{x} \mid \boldsymbol{\theta}) &= \ell(\mathbf{x}, \mathbf{z} \mid \boldsymbol{\theta}) - \log f(\mathbf{z} \mid \mathbf{x}, \boldsymbol{\theta})\end{aligned}$$

Taking expectations of both sides over the distribution of $\mathbf{z} \mid \boldsymbol{\theta}^{(t)}$ gives

$$\ell(\mathbf{x} \mid \boldsymbol{\theta}) = Q(\boldsymbol{\theta} \mid \boldsymbol{\theta}^{(t)}) - H(\boldsymbol{\theta} \mid \boldsymbol{\theta}^{(t)})$$

where

$$\begin{aligned}Q(\boldsymbol{\theta} \mid \boldsymbol{\theta}^{(t)}) &= \int \ell(\mathbf{x}, \mathbf{z} \mid \boldsymbol{\theta}) f(\mathbf{z} \mid \mathbf{x}, \boldsymbol{\theta}^{(t)}) \, d\mathbf{z} \\ H(\boldsymbol{\theta} \mid \boldsymbol{\theta}^{(t)}) &= \int \log f(\mathbf{z} \mid \mathbf{x}, \boldsymbol{\theta}) f(\mathbf{z} \mid \mathbf{x}, \boldsymbol{\theta}^{(t)}) \, d\mathbf{z}\end{aligned}$$

Now, the difference between $\ell(\mathbf{x} \mid \boldsymbol{\theta}^{(t)})$ and $\ell(\mathbf{x} \mid \boldsymbol{\theta}^{(t+1)})$ is given by

$$\begin{aligned}\ell(\mathbf{x} \mid \boldsymbol{\theta}^{(t+1)}) - \ell(\mathbf{x} \mid \boldsymbol{\theta}^{(t)}) &= Q(\boldsymbol{\theta}^{(t+1)} \mid \boldsymbol{\theta}^{(t)}) - H(\boldsymbol{\theta}^{(t+1)} \mid \boldsymbol{\theta}^{(t)}) \\ &\quad - Q(\boldsymbol{\theta}^{(t)} \mid \boldsymbol{\theta}^{(t)}) + H(\boldsymbol{\theta}^{(t)} \mid \boldsymbol{\theta}^{(t)}) \\ &= \left\{ Q(\boldsymbol{\theta}^{(t+1)} \mid \boldsymbol{\theta}^{(t)}) - Q(\boldsymbol{\theta}^{(t)} \mid \boldsymbol{\theta}^{(t)}) \right\} \\ &\quad + \left\{ H(\boldsymbol{\theta}^{(t)} \mid \boldsymbol{\theta}^{(t)}) - H(\boldsymbol{\theta}^{(t+1)} \mid \boldsymbol{\theta}^{(t)}) \right\}\end{aligned}$$

Now:

- The EM algorithm maximises Q , and so $Q(\boldsymbol{\theta}^{(t+1)} \mid \boldsymbol{\theta}^{(t)}) - Q(\boldsymbol{\theta}^{(t)} \mid \boldsymbol{\theta}^{(t)}) > 0$.
- Jensen’s Inequality gives that $H(\boldsymbol{\theta} \mid \boldsymbol{\theta}^{(t)}) \leq H(\boldsymbol{\theta}^{(t)} \mid \boldsymbol{\theta}^{(t)})$.

Overall, therefore, the likelihood increases. ■

Note: The EM algorithm is therefore a “hill-climbing” technique. It is guaranteed to find local maxima. Global maxima can be sought by running the algorithm with many different starting values.

EXAMPLE: Suppose we have a series of data x_1, \dots, x_n from some mixture distribution

$$f(x) = \alpha f_1(x) + (1 - \alpha) f_2(x)$$

For example, UK heights where f_1 corresponds to men, f_2 corresponds to women, and $f_i(x) = N(\mu_i, \sigma^2)$.

The full likelihood function for these data is

$$L(\mathbf{x} | \boldsymbol{\theta}) = \prod_{i=1}^n \left\{ \alpha f_1(x_i) + (1 - \alpha) f_2(x_i) \right\}$$

Maximising this, however, is extremely nasty.

Instead, let

$$z_i = \begin{cases} 1 & \text{if } x_i \sim f_1(x) \\ 0 & \text{otherwise} \end{cases}$$

The likelihood conditional on this variable is then

$$L(\mathbf{x} | \boldsymbol{\theta}, \mathbf{z}) = \prod_{i=1}^n \left\{ \alpha f_1(x_i) \right\}^{z_i} \left\{ (1 - \alpha) f_2(x_i) \right\}^{1-z_i}$$

This is much simpler to maximise.

Let's first consider the E -step. z_i can only be equal to 0 or 1, so

$$\begin{aligned} \mathbb{E}(z_i | \mathbf{x}, \boldsymbol{\theta}) &= 1 \times \mathbb{P}(z_i = 1 | \mathbf{x}, \boldsymbol{\theta}) + 0 \times \mathbb{P}(z_i = 0 | \mathbf{x}, \boldsymbol{\theta}) \\ &= \mathbb{P}(z_i = 1 | \mathbf{x}, \boldsymbol{\theta}) \\ &= \frac{\alpha f_1(x_i)}{\alpha f_1(x_i) + (1 - \alpha) f_2(x_i)} = \hat{z}_i \end{aligned}$$

(This is evident from the form of the likelihood, which makes it clear that z is a Bernoulli random variable.)

Now the M -step. We have that $Q(\boldsymbol{\theta} | \boldsymbol{\theta}^{(t)}) = \ell(\mathbf{x}, \hat{\mathbf{z}} | \boldsymbol{\theta})$

$$\ell(\mathbf{x} | \boldsymbol{\theta}, \mathbf{z}) = \sum_{i=1}^n \left[z_i \log \left\{ \alpha f_1(x_i) \right\} + (1 - z_i) \log \left\{ (1 - \alpha) f_2(x_i) \right\} \right]$$

and

$$\begin{aligned} \frac{\partial \ell}{\partial \alpha} &= \sum_{i=1}^n \left[\frac{z_i}{\alpha} + \frac{-(1 - z_i)}{(1 - \alpha)} \right] = \sum_{i=1}^n \left[\frac{z_i - \alpha}{\alpha(1 - \alpha)} \right] = 0 \\ &\Rightarrow \hat{\alpha} = \frac{\sum_{i=1}^n z_i}{n} \end{aligned}$$

and

$$\begin{aligned} \frac{\partial \ell}{\partial \mu_1} &= -\frac{\partial}{\partial \mu_1} \sum_{i=1}^n z_i \frac{1}{2} \left(\frac{x_i - \mu_1}{\sigma} \right)^2 = \sum_{i=1}^n z_i \frac{x_i - \mu_1}{\sigma^2} = 0 \\ &\Rightarrow \hat{\mu}_1 = \frac{\sum_{i=1}^n z_i x_i}{\sum_{i=1}^n z_i}, \quad \hat{\mu}_2 = \frac{\sum_{i=1}^n (1 - z_i) x_i}{\sum_{i=1}^n (1 - z_i)} \end{aligned}$$

and finally

$$\begin{aligned}\frac{\partial \ell}{\partial \sigma^2} &= \sum_{i=1}^n \left[z_1 \left\{ -\frac{1}{2\sigma^2} + \frac{(x - \mu_1)^2}{2(\sigma^2)^2} \right\} + (1 - z_1) \left\{ -\frac{1}{2\sigma^2} + \frac{(x - \mu_2)^2}{2(\sigma^2)^2} \right\} \right] \\ &= \sum_{i=1}^n \left[\frac{z_1(x - \mu_1)^2}{2(\sigma^2)^2} - \frac{1}{2\sigma^2} + \frac{(x - \mu_2)^2}{2(\sigma^2)^2} - \frac{z_1(x - \mu_2)^2}{2(\sigma^2)^2} \right] = 0\end{aligned}$$

and so

$$\begin{aligned}\sum_{i=1}^n \left[\frac{z_1(x - \mu_1)^2}{2(\hat{\sigma}^2)^2} - \frac{1}{2\hat{\sigma}^2} + \frac{(x - \mu_2)^2}{2(\hat{\sigma}^2)^2} - \frac{z_1(x - \mu_2)^2}{2(\hat{\sigma}^2)^2} \right] &= 0 \\ \sum_{i=1}^n \left[\frac{z_1(x - \mu_1)^2}{2(\hat{\sigma}^2)^2} + \frac{(x - \mu_2)^2}{2(\hat{\sigma}^2)^2} - \frac{z_1(x - \mu_2)^2}{2(\hat{\sigma}^2)^2} \right] &= \frac{n}{2\hat{\sigma}^2} \\ \hat{\sigma}^2 &= \frac{1}{n} \sum_{i=1}^n [z_1(x - \mu_1)^2 + (1 - z_1)(x - \mu_2)^2]\end{aligned}$$

Thus, our algorithm simply involves finding z in each case (at the E -step) and then finding values for the other parameters (at the M -step).