

# STAT 5361: Statistical Computing

Jun Yan

Department of Statistics  
University of Connecticut

# Outline I

## 1 Non-Stochastic Optimization

- Some Background
- Univariate Problems
- Multivariate Problems

## 2 Combinatorial Optimization

- Local Search
- Simulated Annealing
- Genetic Algorithm

## 3 EM and MM Algorithms

- EM Algorithm
- Majorization-Minimization Algorithm
  - Example: Penalized AFT model with induced smoothing

## 4 Simulation Basics

- Inverse Transform Method
- Rejection Sampling
- Transformation
- Variants of Rejection Sampling

# Outline II

- Sampling importance resampling
- Sampling using Markov Chains

## 5 Simulation of Stochastic Processes

- Simulating Brownian motions
- Gaussian short rate models
- Processes with Jumps

## 6 Monte Carlo Integration

- Basics of Monte Carlo Integration
- State-Space model
- Efficiency of MC

# Non-stochastic Optimization

What we will learn:

- Some basics about statistical inference
- Univariate problems
  - ▶ Newton's method
  - ▶ Fisher scoring
  - ▶ Secant method
  - ▶ Fixed-point method
  - ▶ Connections of the above
- Multivariate problems
  - ▶ Newton's method
  - ▶ Newton-like methods

## Background: likelihood inference

- Let  $\mathbf{x}_1, \dots, \mathbf{x}_n$  be an i.i.d. sample from

$$f(\mathbf{x} | \boldsymbol{\theta}^*),$$

with the true parameter value  $\boldsymbol{\theta}^*$  being unknown.

- The likelihood function is

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

- The *maximum likelihood estimator* (MLE) of the parameter value is the maximizer of  $L(\boldsymbol{\theta})$ . MLE has the invariate property.
- Usually it is easier to work with the *log likelihood function*

$$l(\boldsymbol{\theta}) = \log L(\boldsymbol{\theta}).$$

- Typically, maximization of  $l(\boldsymbol{\theta})$  is done by solving

$$l'(\boldsymbol{\theta}) = 0.$$

- $l'(\boldsymbol{\theta})$  is called the *score function*.
  - For each possible parameter value  $\boldsymbol{\theta}$ ,  $l(\boldsymbol{\theta})$  is a random variable, because it depends on the observed values of  $\mathbf{x}_1, \dots, \mathbf{x}_n$ .

- For any  $\boldsymbol{\theta}$ ,

$$\begin{aligned}\mathbb{E}_{\boldsymbol{\theta}} \{l'(\boldsymbol{\theta})\} &= 0, \\ \mathbb{E}_{\boldsymbol{\theta}} \{l'(\boldsymbol{\theta})l'(\boldsymbol{\theta})^T\} &= -\mathbb{E}_{\boldsymbol{\theta}} \{l''(\boldsymbol{\theta})\}.\end{aligned}$$

where  $\mathbb{E}_{\boldsymbol{\theta}}$  is the expectation with respect to  $f(\mathbf{x} | \boldsymbol{\theta})$ . The equality holds true under mild regularity conditions.



$$I(\boldsymbol{\theta}) = \mathbb{E}_{\boldsymbol{\theta}} \{l'(\boldsymbol{\theta})l'(\boldsymbol{\theta})^T\}$$

is known as the *Fisher information*.

- The importance of the Fisher information  $I(\theta)$  is that it sets the limit on how accurate an unbiased estimate of  $\theta$  can be.
- As  $n \rightarrow \infty$ , the asymptotic distribution of  $\sqrt{n}(\hat{\theta}_{\text{MLE}} - \theta^*)$  is  $N_p(\mathbf{0}, nI(\theta^*)^{-1})$ . Since  $\theta^*$  is unknown, the asymptotic covariance matrix  $I(\theta^*)^{-1}$  needs to be estimated.
  - ▶ If  $\dim(\theta) = 1$ ,  $I(\theta)$  is a nonnegative number.
  - ▶ If  $\dim(\theta) > 1$ ,  $I(\theta)$  is a nonnegative definite matrix.
- The *observed Fisher information* is

$$-l''(\theta).$$

The expected Fisher information  $I(\theta)$  may not be easily computed. The observed Fisher information is a good approximation to  $I(\theta)$  that improves as  $n$  gets bigger.

Besides MLEs, there are other likelihoods for parameter estimation. Suppose  $\theta$  has two parts:  $\theta = (\phi, \mu)$  and we are only interested in  $\phi$ . The *profile likelihood* for  $\phi$  is

$$L(\phi) := \max_{\mu} L(\phi, \mu).$$

- $\mu$  is the nuisance parameter.
- Need to maximize  $L(\phi, \mu)$  for every fixed  $\phi$ . Also need to maximize  $L(\phi)$ .



## A general method

Suppose  $g(\mathbf{x})$  is a differentiable function, where

$$\mathbf{x} = (x_1, \dots, x_n).$$

To find its maximum (or minimum), one method is to solve the equation

$$g'(\mathbf{x}) = 0$$

$$\text{where } g'(\mathbf{x}) = \left( \frac{\partial g}{\partial x_1}, \dots, \frac{\partial g}{\partial x_n} \right)^T.$$

Then the maximization is equivalent to solving

$$f(\mathbf{x}) = 0,$$

where  $f = g'$ .

For maximum likelihood estimation,  $g$  is the log likelihood function  $l$ , and  $\mathbf{x}$  is the corresponding parameter vector  $\boldsymbol{\theta}$ .

## Univariate Problems

## Optimization: Univariate case

- Goal: optimize a real-valued function  $g$  with respect to its argument, a  $p$  dimensional vector  $\mathbf{x}$ . We will first consider the case  $p = 1$ .
- We will limit consideration mainly to smooth and differentiable functions.
- Root finding methods. Solving unconstrained nonlinear equations.
- Iterative algorithms: starting value, updating equation, stopping rule/convergence criterion.

Using bisection method to maximize

$$g(x) = \frac{\log x}{1 + x}.$$

or to solve

$$f(x) = g'(x) = \frac{1 + \frac{1}{x} - \log x}{(1 + x)^2} = 0.$$

# Newton's method

This is a fast approach to finding roots of a differentiable function  $f(x)$ . First, set initial value  $x_0$ . Then for  $t = 0, 1, \dots$ , compute

$$x_{t+1} = x_t + h_t, \text{ with } h_t = -\frac{f(x_t)}{f'(x_t)}.$$

Continue the iterations until  $x_t$  converges.

- Also known as Newton-Raphson iteration.
- Need to specify  $x_0$ .
- If  $f(x) = 0$  has multiple solutions, the end result depends on  $x_0$ .

Newton's method require computing the derivative of a function. Algorithms are available to find root(s) of a univariate function without having to compute its derivative. For example, in R, one can use `uniroot`. Newton's method can be applied to optimize  $g$  by applying to

$$f = g'.$$

- Both  $g'$  (*gradient*) and  $g''$  (*Hessian*) are needed.
- Many variants of Newton's method avoid the computation of the Hessian, which can be difficult especially for multivariate functions.

To maximize

$$g(x) = \frac{\log x}{1+x},$$

first find

$$f(x) = g'(x) = \frac{1 + \frac{1}{x} - \log x}{(1+x)^2},$$

$$f'(x) = g''(x) = -\frac{3 + 4/x + 1/x^2 - 2 \log x}{(1+x)^3}.$$

So in the Newton's method,

$$h_t = \frac{(x_t + 1)(1 + 1/x_t - \log x_t)}{3 + 4/x_t + 1/x_t^2 - 2 \log x_t}.$$

Note that to solve  $f(x) = 0$ , one can instead solve  $1 + 1/x - \log x = 0$ . Treat this as a new  $f$  function. Then in the Newton's method,

$$h_t = x_t - \frac{x_t^2 \log x_t}{1 + x_t} \implies x_{t+1} = 2x_t - \frac{x_t^2 \log x_t}{1 + x_t}.$$

To maximize the log likelihood  $l(\theta)$ , Newton's method computes

$$\theta_{t+1} = \theta_t - \frac{l'(\theta_t)}{l''(\theta_t)}$$

Example: consider  $x_1, \dots, x_n \sim \text{i.i.d. } N(\mu, \sigma^2)$ .

Example: consider the model on shift

$$p(x | \theta) = p(x - \theta).$$

Given observations  $x_1, \dots, x_n$  i.i.d.  $\sim p(x | \theta)$

$$l(\theta) = \sum_{i=1}^n \log p(x_i - \theta)$$

$$l'(\theta) = - \sum_{i=1}^n \frac{p'(x_i - \theta)}{p(x_i - \theta)}$$

$$l''(\theta) = \sum_{i=1}^n \frac{p''(x_i - \theta)}{p(x_i - \theta)} - \sum_{i=1}^n \left\{ \frac{p'(x_i - \theta)}{p(x_i - \theta)} \right\}^2.$$

Note that we need to update  $\theta$  in the iterations, not  $x_1, \dots, x_n$ .



In R, to *minimize* a function, one can use

```
z = nlminb(x0, g, gr.g, hess.g)
```

here  $x_0$  is the initial value,  $g$  is the function being minimized,  $gr.g$  its gradient and  $hess.g$  its Hessian. (Unconstrained and box-constrained optimization using PORT routines).

In the above function, the derivatives  $g'$  and  $g''$  have to be analytically calculated. One can also use

```
z = nlminb(x0, g, gr.g)
```

without inputting the analytic expression of  $g''$ , or, even simpler,

```
z = nlminb(x0, g)
```

without inputting the analytic expressions of either derivatives. In these cases, numerical approximations of derivatives will be computed during the iterations.

Other functions for optimization in R include `optim`, `optimize`, `nlm` and `constrOptim`.

For profile likelihood

$$L(\phi) = \max_{\mu} L(\phi, \mu)$$

we need to optimize  $L(\phi, \mu)$  for each given  $\phi$ .

In R, in order to get  $\min_x g(x, y)$  for each fixed  $y$ , one can do the following. First, define  $g(x, y)$ , `gr.g(x, y)`, etc. Then call, say,

```
nlminb(x0, g, gr.g, hess.g, y=1), or  
nlminb(x0, g, gr.g, y=.3), or  
nlminb(x0, g, y=-1)
```

If one only wants minimization for  $x \in [a, b]$ , use,

```
nlminb(x0, ..., lower=a, upper=b)
```

## Fisher scoring

This is a variant of Newton's method specific for MLE. Recall  $-l''(\theta)$  is the observed Fisher information at  $\theta$ . To maximize  $l(\theta)$ , an alternative is to replace  $-l''(\theta_t)$  with  $I(\theta_t)$  to get

$$\theta_{t+1} = \theta_t + \frac{l'(\theta_t)}{I(\theta_t)}.$$

To *minimize*  $-l(\theta)$ , one still can use

$$z = \text{nlminb}(x0, f, \text{grf}, \text{fs})$$

where  $f$  is  $-l(\theta)$ ,  $\text{grf}$  is  $-l'(\theta)$ , and  $\text{fs}$  is  $I(\theta)$  instead of  $-l''(\theta)$ .

Generally, use Fisher scoring in the beginning to make rapid improvements, and Newton's method for refinement near the end.

Continuing the example on  $p(x | \theta) = p(x - \theta)$ , to use Fisher scoring, we need to compute

$$I(\theta) = -\mathbb{E}_{\theta}[l''(\theta)].$$

We already know

$$l''(\theta) = \sum_{i=1}^n \frac{p''(x_i - \theta)}{p(x_i - \theta)} - \sum_{i=1}^n \left\{ \frac{p'(x_i - \theta)}{p(x_i - \theta)} \right\}^2.$$

Therefore

$$I(\theta) = -n\mathbb{E}_{\theta} \left[ \frac{p''(X - \theta)}{p(X - \theta)} - \left\{ \frac{p'(X - \theta)}{p(X - \theta)} \right\}^2 \right].$$

Since under parameter  $\theta$ ,  $X$  has density  $p(x - \theta)$ , the last  $\mathbb{E}_\theta[\dots]$  equals

$$\begin{aligned} & \int \left[ \frac{p''(x - \theta)}{p(x - \theta)} - \left\{ \frac{p'(x - \theta)}{p(x - \theta)} \right\}^2 \right] p(x - \theta) \, dx \\ &= \int p''(x - \theta) \, dx - \int \frac{[p'(x - \theta)]^2}{p(x - \theta)} \, dx \\ &= \frac{d^2}{d\theta^2} \int p(x - \theta) \, dx - \int \frac{\{p'(x)\}^2}{p(x)} \, dx \\ &= \frac{d^2}{d\theta^2} 1 - \int \frac{\{p'(x)\}^2}{p(x)} \, dx = - \int \frac{\{p'(x)\}^2}{p(x)} \, dx. \end{aligned}$$

Therefore,

$$I(\theta) = n \int \frac{\{p'(x)\}^2}{p(x)} \, dx.$$

So, in this case, Fisher information is a constant.

# Secant method

Approximating  $f'(x_t)$  by

$$\frac{f(x_t) - f(x_{t-1})}{x_t - x_{t-1}},$$

the Newton's method turns into the secant method

$$x_{t+1} = x_t - \frac{f(x_t)(x_t - x_{t-1})}{f(x_t) - f(x_{t-1})}.$$

- Need to specify  $x_0$  and  $x_1$ .

# Fixed point iteration

Compute

$$x_{t+1} = x_t + \alpha f(x_t), \quad t = 0, 1, \dots,$$

where  $\alpha \neq 0$  is a tuning parameter so that

$$|1 + \alpha f'(x)| \leq \lambda < 1, \quad \text{all } x$$

or more generally,

$$|x - y + \alpha[f(x) - f(y)]| \leq \lambda|x - y|, \quad \text{any } x, y$$

with  $0 \leq \lambda < 1$  a constant, i.e.,  $F(x) = x + \alpha f(x)$  is a *contraction*.

If such  $\alpha$  exists, the speed of convergence depends on  $\lambda$ . The smaller  $\lambda$  is, the faster the convergence.

## Some Details on Fixed point iteration

Definition: A fixed point of a function is a point whose evaluation by that function equals to itself, i.e.,  $x = G(x)$ .

Fixed point iteration: the natural way of hunting for a fixed point is to use  $x_{t+1} = G(x_t)$ .

Definition: A function  $G$  is contractive on  $[a, b]$  if

- (1).  $G(x) \in [a, b]$  whenever  $x \in [a, b]$ ,
- (2).  $|G(x_1) - G(x_2)| \leq \lambda |x_1 - x_2|$  for all  $x_1, x_2 \in [a, b]$  and some  $\lambda \in [0, 1)$ .

Theorem: if  $G$  is contractive on  $[a, b]$ , then there is a unique fixed point  $x^* \in [a, b]$ , and the fixed point iteration convergence to it when starting inside the interval.

Convergence:

$|x_{t+1} - x_t| = |G(x_t) - G(x_{t-1})| \leq \lambda |x_t - x_{t-1}| \leq \lambda^t |x_1 - x_0| \rightarrow 0$ , as  $t \rightarrow \infty$ . It follows that  $\{x_t\}$  convergent to a limit  $x^*$ .



# Connection between fixed-point method and Newton methods

Root-finding problems using fixed point iteration: for solving  $f(x) = 0$ , we can simply let  $G(x) = x + \alpha f(x)$ , where  $\alpha \neq 0$  is a constant.

Required Lipschitz condition:  $|x - y + \alpha[f(x) - f(y)]| \leq \lambda|x - y|$ , for some  $\lambda \in [0, 1)$  and for all  $x, y \in [a, b]$ . This holds if  $|G'(x)| \leq \lambda$  for some  $\lambda \in [0, 1)$  and for all  $x \in [a, b]$ , i.e.,  $|1 + \alpha f'(x)| \leq \lambda$ . (use mean value theorem.)

Newton methods:  $G(x) = x - f(x)/f'(x)$ . So it is as if  $\alpha_t$  is chosen adaptively as  $\alpha_t = -1/f'(x_t)$ . This leads to a faster convergence order (quadratic).

# Convergence Order

Define  $\varepsilon_t = x_t - x^*$ . A method has convergence order  $\beta$  if

$$\lim_{t \rightarrow \infty} \varepsilon_t = 0, \quad \lim_{t \rightarrow \infty} \frac{|\varepsilon_{t+1}|}{|\varepsilon_t|^\beta} = c,$$

for some constant  $c \neq 0$  and  $\beta > 0$ .

- For Newton's method,  $\beta = 2$ . (If it converges.)
- For Secant method,  $\beta \approx 1.62$ .
- For fixed point iteration,  $\beta = 1$ .

## Convergence Order

For Newton's method: suppose  $f$  has two continuous derivatives and  $f'(x^*) \neq 0$ . There then exists a neighborhood of  $x^*$  within which  $f'(x) \neq 0$  for all  $x$ . By Taylor expansion,

$$0 = f(x^*) = f(x_t) + f'(x_t)(x^* - x_t) + \frac{1}{2}f''(q)(x^* - x_t)^2,$$

for some  $q$  between  $x^*$  and  $x_t$ . Rearranging terms, we find that

$$\frac{\varepsilon_{t+1}}{\varepsilon_t^2} = \frac{f''(q)}{2f'(x_t)} \rightarrow \frac{f''(x^*)}{2f'(x^*)}.$$

# Convergence Order

For fixed point iteration: let  $G$  be a continuous function on the closed interval  $[a, b]$  with  $G : [a, b] \rightarrow [a, b]$  and suppose that  $G'$  is continuous on the open interval  $(a, b)$  with  $|G'(x)| \leq k < 1$  for all  $x \in (a, b)$ . If  $G'(x^*) \neq 0$ , then for any  $x_0 \in [a, b]$ , the fixed point iteration converges linearly to the fixed point  $x^*$ . This is because

$$|x_{t+1} - x^*| = |G(x_t) - G(x^*)| = |G'(q)||x_t - x^*|,$$

for some  $q$  between  $x_t$  and  $x^*$ . This implies that

$$\frac{|\varepsilon_{t+1}|}{|\varepsilon_t|} \rightarrow |G'(x^*)|.$$

# Stopping Rules

- Absolute convergence criterion:

$$\|x_{t+1} - x_t\| < \epsilon,$$

where  $\epsilon$  is a chosen tolerance.

- Relative convergence criterion

$$\frac{\|x_{t+1} - x_t\|}{\|x_t\|} < \epsilon.$$

If  $x_t$  close to zero,

$$\frac{\|x_{t+1} - x_t\|}{\|x_t\| + \epsilon} < \epsilon.$$

- Many other rules depending on the algorithm.

## Multivariate Problems

# Newton's method

Let  $g$  now be a function in  $\mathbf{x} = (x_1, \dots, x_p)^T$ .

The generalization is straightforward: to maximize  $g(\mathbf{x})$ , set

$$\mathbf{x}_{t+1} = \mathbf{x}_t - [g''(\mathbf{x}_t)]^{-1} g'(\mathbf{x}_t).$$

- $g''(\mathbf{x})$  is a  $p \times p$  matrix with the  $(i, j)$ th entry equal to

$$\frac{\partial^2 g(\mathbf{x})}{\partial x_i \partial x_j};$$

- $g'(\mathbf{x})$  is a  $p \times 1$  vector with the  $i$ th entry equal to

$$\frac{\partial g(\mathbf{x})}{\partial x_i}.$$

## Newton-like methods

For MLE, the generalization is straightforward. Simply use

$$\boldsymbol{\theta}_{t+1} = \boldsymbol{\theta}_t + I(\boldsymbol{\theta}_t)^{-1} l'(\boldsymbol{\theta}_t).$$

These methods in each iteration approximate  $g''(\mathbf{x}_t)$  by some  $p \times p$  matrix  $\mathbf{M}_t = \mathbf{M}_t(\mathbf{x}_t)$  to get

$$\mathbf{x}_{t+1} = \mathbf{x}_t - \mathbf{M}_t^{-1} g'(\mathbf{x}_t). \quad (1)$$

Of course,  $\mathbf{M}_t$  should be easier to compute than  $g''(\mathbf{x}_t)$ .

Fisher scoring is a Newton-like method, because it uses

$$\mathbf{M}_t = -I(\boldsymbol{\theta}_t)$$

in place of  $-l''(\boldsymbol{\theta}_t)$ .



# Steepest ascent methods

In (1), set

$$\mathbf{M}_t = -\alpha_t^{-1} \mathbf{I}_p,$$

so that

$$\mathbf{x}_{t+1} = \mathbf{x}_t + \alpha_t g'(\mathbf{x}_t),$$

where  $\alpha_t > 0$  is the step size at  $t$  which can shrink to ensure ascent. If at step  $t$ , the original step turns out to be downhill, the updating can backtrack by halving  $\alpha_t$ .

## Discrete Newton and Fixed-point methods

In fixed-point methods, usually  $\mathbf{M}_t$  is fixed to be  $\mathbf{M}$ , e.g.,  $g''(\mathbf{x}_0)$ . This amounts to applying univariate scaled fixed-point iteration to each component.

In discrete Newton methods,  $g''(\mathbf{x}_t)$  is approximated as follows. Compute matrix  $\mathbf{M}_t$  with the  $(i, j)^{\text{th}}$  entry equal to

$$\mathbf{M}_t(i, j) = \frac{g'_i(\mathbf{x}_t + h_t(i, j)\mathbf{e}_j) - g'_i(\mathbf{x}_t)}{h_t(i, j)}$$

for some constant  $h_t(i, j) \neq 0$ , where

$$g'_i(\mathbf{x}) = \frac{\partial g(\mathbf{x})}{\partial x_i}.$$

If  $h_t(i, j)$  is small, then

$$\mathbf{M}_t(i, j) \approx \frac{\partial^2 g(\mathbf{x}_t)}{\partial x_i \partial x_j}$$

and so  $\mathbf{M}_t$  approximates  $g''(\mathbf{x}_t)$ .

However, since  $g''(\mathbf{x}_t)$  is symmetric, instead of using  $\mathbf{M}_t$  directly, use the symmetric

$$(\mathbf{M}_t + \mathbf{M}_t^T)/2$$

as the approximation of  $g''(\mathbf{x}_t)$ .

- No need to calculate second order derivatives
- Inefficient: all  $p^2$  entries have to be updated each time.

# Quasi-Newton methods

These methods aim to achieve the following goals.

- Each step satisfies the secant condition

$$g'(\mathbf{x}_{t+1}) - g'(\mathbf{x}_t) = \mathbf{M}_{t+1}(\mathbf{x}_{t+1} - \mathbf{x}_t).$$

- No need to compute second order derivatives.
- Maintain symmetry of  $\mathbf{M}_t$ .
- Aim to update  $\mathbf{M}_t$  efficiently.

There is a unique rank-one method that satisfies the secant condition and maintains the symmetry of  $M_t$ : after getting

$$\mathbf{x}_{t+1} = \mathbf{x}_t - M_t^{-1} g'(\mathbf{x}_t)$$

compute

$$\begin{aligned} \mathbf{z}_t &= \mathbf{x}_{t+1} - \mathbf{x}_t, & \mathbf{y}_t &= g'(\mathbf{x}_{t+1}) - g'(\mathbf{x}_t), \\ \mathbf{v}_t &= \mathbf{y}_t - M_t \mathbf{z}_t, & c_t &= \frac{1}{\mathbf{v}_t^T \mathbf{z}_t}. \end{aligned}$$

Then update

$$M_{t+1} = M_t + c_t \mathbf{v}_t \mathbf{v}_t^T.$$

Note  $M_{t+1} - M_t$  is of rank one. We can verify the secant condition is satisfied by multiplying both sides by  $\mathbf{z}_t$ .

There are several rank-two method satisfying the secant condition and the symmetry requirement. The Broyden class updates  $\mathbf{M}_t$  as follows. After getting  $\mathbf{x}_{t+1}$  and  $\mathbf{z}_t$ ,  $\mathbf{y}_t$  as above, compute

$$\mathbf{d}_t = \frac{\mathbf{y}_t}{\mathbf{z}_t^T \mathbf{y}_t} - \frac{\mathbf{M}_t \mathbf{z}_t}{\mathbf{z}_t^T \mathbf{M}_t \mathbf{z}_t}.$$

Then update

$$\mathbf{M}_{t+1} = \mathbf{M}_t - \frac{\mathbf{M}_t \mathbf{z}_t (\mathbf{M}_t \mathbf{z}_t)^T}{\mathbf{z}_t^T \mathbf{M}_t \mathbf{z}_t} + \frac{\mathbf{y}_t \mathbf{y}_t^T}{\mathbf{z}_t^T \mathbf{y}_t} + \delta_t (\mathbf{z}_t^T \mathbf{M}_t \mathbf{z}_t) \mathbf{d}_t \mathbf{d}_t^T,$$

where  $\delta_t$  is a parameter.

A popular method to solve unconstrained nonlinear optimization problems is the Broyden-Fletcher-Goldfarb-Shanno (BFGS) method, with  $\delta_t \equiv 0$ .

In R, `optim` can be called to minimize a function, using the BFGS method or its variant, the “L-BFGS-B” method. “L” stands for limited memory, and “B” stands for box constraint.

# Approximating Hessian for MLE

For MLE, the Hessian  $l''(\hat{\theta}_{\text{MLE}})$  is critical because it provides estimates of standard error and covariance.

- Quasi-Newton methods may provide poor approximation because it is based on the idea of using poor approximation of the Hessian to find the root for  $l'(\theta) = 0$ .
- Use the discrete multivariate Newton method with

$$M_t(i, j) = \frac{l'_i(\theta_t + h_{ij}e_j) - l'_i(\theta_t - h_{ij}e_j)}{2h_{ij}}.$$

## Gauss-Newton method

Example: nonlinear regression. In many cases, we want to maximize

$$g(\boldsymbol{\theta}) = -\sum_{i=1}^n (y_i - f_i(\boldsymbol{\theta}))^2,$$

where each  $f_i(\boldsymbol{\theta})$  is differentiable. Recall that for linear regression:

$$y_i = \mathbf{x}_i^T \boldsymbol{\theta} + \varepsilon, \quad i = 1, \dots, n$$

the LS estimator of  $\boldsymbol{\theta}$  maximizes  $g(\boldsymbol{\theta})$  with  $f_i(\boldsymbol{\theta}) = \mathbf{x}_i^T \boldsymbol{\theta}$  and equals

$$\hat{\boldsymbol{\theta}} = (\mathbf{X}^T \mathbf{X})^{-1} \mathbf{X}^T \mathbf{y}$$

where

$$\mathbf{X} = \begin{pmatrix} \mathbf{x}_1^T \\ \vdots \\ \mathbf{x}_n^T \end{pmatrix}, \quad \mathbf{y} = \begin{pmatrix} y_1 \\ \vdots \\ y_n \end{pmatrix}.$$



The Gauss-Newton method applies a similar idea to the nonlinear case. Let  $\theta^*$  be the (unknown) maximizer of  $g(\theta)$ . Given any candidate  $\theta$ , the function

$$h(u) = - \sum_{i=1}^n (y_i - f_i(\theta + u))^2.$$

is maximized by  $\beta = \theta^* - \theta$ . Of course,  $\beta$  is unknown. However, if  $\theta$  is near  $\theta^*$ , then  $\beta \approx 0$ , so by Taylor expansion of  $h(u)$ , it may be close to the maximizer of

$$- \sum_{i=1}^n (y_i - f_i(\theta) - f'_i(\theta)^T u)^2.$$

Treat  $y_i - f_i(\boldsymbol{\theta})$  the same way as  $y_i$  in the linear regression, and  $f'_i(\boldsymbol{\theta})$  the same way as  $\mathbf{x}_i$ , then

$$\boldsymbol{\theta}^* - \boldsymbol{\theta} \approx (\mathbf{A}^T \mathbf{A})^{-1} \mathbf{A}^T \mathbf{z}$$

where

$$\mathbf{z} = \mathbf{z}(\boldsymbol{\theta}) = \begin{pmatrix} y_1 - f_1(\boldsymbol{\theta}) \\ \vdots \\ y_n - f_n(\boldsymbol{\theta}) \end{pmatrix}, \quad \mathbf{A} = \mathbf{A}(\boldsymbol{\theta}) = \begin{pmatrix} f'_1(\boldsymbol{\theta})^T \\ \vdots \\ f'_n(\boldsymbol{\theta})^T \end{pmatrix}.$$

The update rule thus is

$$\boldsymbol{\theta}_{t+1} = \boldsymbol{\theta}_t + (\mathbf{A}_t^T \mathbf{A}_t)^{-1} \mathbf{A}_t^T \mathbf{z}_t,$$

where  $\mathbf{z}_t = \mathbf{z}(\boldsymbol{\theta}_t)$  and  $\mathbf{A}_t = \mathbf{A}(\boldsymbol{\theta}_t)$ .

In R, the Gauss-Newton method is the default method of the function `nls`.

# Practical Issues

- Initial value
- Convergence criteria: based on a distance measures for vectors
- Multiple local optima
  - ▶ Multiple initial values
  - ▶ Compare the objective function value at convergences
- Trade-off between efficiency and robustness

# Combinatorial Optimization

## Motivation:

- There are hard optimization problems for which most methods including what we have discussed so far are useless.
- These problems are usually combinatorial in nature. Maximization requires a discrete search of a very large space.
- Problem: maximize  $f(\theta)$  with respect to  $\theta = (\theta_1, \dots, \theta_p)$ , where  $\theta \in \Theta$  and  $\Theta$  consists of  $N$  elements.
  - ▶ Each  $\theta \in \Theta$  is called a candidate solution.
  - ▶ Usually  $N$  is a very large number depending on problem size  $p$ .
  - ▶ The difficulty of a particular size- $p$  problem can be characterized by the number of operations required to solve it in the worst case scenario using the best known algorithm.
  - ▶ Suppose a problem is  $\mathcal{O}(p!)$ . If it requires 1 minute to solve for  $p = 20$ , it would take 12.1 years for  $p = 25$  and 207 million years for  $p = 30$ .

What can we do:

- We have to take some sacrifices: we will abandon the global algorithms and focus on algorithms that can find a good local solution.
- Heuristic strategies.

What we will learn:

- Local search methods
- Simulated annealing
- Genetic algorithms

## Motivating Example: Subset regression

Suppose  $x_1, \dots, x_p$  are predictors that can be used to construct a linear model for  $Y$ . Each candidate model is

$$Y = \sum_{j=1}^s \beta_{i_j} x_{i_j} + \varepsilon,$$

where  $1 \leq i_1 < i_2 < \dots < i_s \leq p$ ,  $s \geq 0$ . The Akaike information criterion (AIC) for the candidate model is

$$\text{AIC} = n \log \frac{\text{RSS}}{n} + 2s$$

where  $n$  is the sample size and RSS is the residual sum of squares of model. The best model is the one that minimizes AIC.

To parameterize the model, set  $\theta = (\theta_1, \dots, \theta_p)$ , such that  $\theta_i = 1$  if  $x_i$  is included as a predictor, and 0 otherwise. The set  $\Theta$  of all possible  $\theta$  has  $2^p$  values.

# Local search

First, define a neighborhood  $\mathcal{N}(\theta)$  for each  $\theta$ , so that it

- contains candidate solutions that are “near”  $\theta$ , and
- reduces the number of changes to the current  $\theta$ .

For example, if  $\Theta = \{\theta = (\theta_1, \dots, \theta_p) : \text{each } \theta_i = 0, 1\}$ , one may define  $\mathcal{N}(\theta)$  to be the set of  $\theta'$  which are different from  $\theta$  in at most one coordinate.

After neighborhoods are defined, at iteration  $t$ , choose  $\theta_{t+1}$  from  $\mathcal{N}(\theta_t)$  according to a certain rule. For example,

- steepest ascent:  $\theta_{t+1} = \arg \max_{\theta \in \mathcal{N}(\theta_t)} f(\theta)$ ;
- ascent algorithm:  $\theta_{t+1} \in \mathcal{N}(\theta_t)$  uphill from  $\theta_t$ , i.e.,

$$f(\theta_{t+1}) \geq f(\theta_t).$$

To avoid trapping into a local maximum, two often used variants are

- random-starts local search: repeatedly run an ascent algorithm to termination from a large number of randomly chosen starting points.
- steepest ascent/mildest descent: set to the least unfavorable  $\theta_{t+1} \in \mathcal{N}(\theta_t)$ ;
  - ▶ if  $\theta_t$  is a local maximum  $\theta_{t+1}$  is the one with least decrease;
  - ▶ otherwise,  $\theta_{t+1}$  is the one with the largest increase.



To select a linear model to minimize AIC (or maximize  $-AIC$ ),

- randomly select a set of predictors
- at iteration  $t$ , decrease the AIC by adding a predictor to the set of selected predictors or deleting a predictor that has been selected; continue the iterations until the AIC can not be decreased;
- repeat Steps 1 and 2 many times, and choose the set of predictors at termination that has the lowest AIC.

A salesman must visit each of  $p$  cities exactly once and return to his starting city, using the shortest total travel distance.

A candidate solution is

$$\theta = (i_1, i_2, \dots, i_p, i_1)$$

where  $i_1, \dots, i_p$  is a permutation of  $1, \dots, p$ . The objective function to be minimized is

$$f(\theta) = d(i_1, i_2) + d(i_2, i_3) + \dots + d(i_{p-1}, i_p) + d(i_p, i_1).$$

There are  $(p-1)!/2$  all possible routes, since the point of origin and direction of travel are arbitrary. For a traveling plan to visit 20 cities, that amounts to more than  $6 \times 10^{16}$  possible routes.

One could define  $\mathcal{N}(\theta)$  as the set of sequences that only differ from  $\theta$  at two entries. In other words, each sequence in  $\mathcal{N}(\theta)$  is obtained by exchanging two entries of  $\theta$ .

For example, if

$$\theta = (1, 2, 5, 4, 6, 3, 1)$$

then

$$(1, 2, 3, 4, 6, 5, 1)$$

is a neighbor of  $\theta$ , because it only has 3 and 5 exchanged; while

$$(1, 2, 4, 3, 6, 5, 1)$$

is not a neighbor of  $\theta$ , because it has 5, 4, and 3 rotated.

# Simulated annealing

In most cases, the simulated annealing algorithm can be thought of as a randomized local search algorithm. It uses a “temperature” parameter to control the randomness of the search. The algorithm starts with a high temperature and cools down gradually so that a global optimum may be reached. This is analogous to the annealing of metal or glass.

In the following description, a global *minimum* of  $f$  is being searched. The algorithm is run in stages, such that for the iterations within a stage, the temperature is a constant  $\tau_j$ .

Suppose iteration  $t$  belongs to stage  $j$ .

1. Sample a candidate  $\theta^* \in \mathcal{N}(\theta)$  according to a *proposal density*  $g_t(\theta | \theta_t)$ ; for different  $t$ , the proposal density  $g_t$  can be different.
2. Let  $\Delta = f(\theta^*) - f(\theta_t)$ .
  - a) If  $\Delta \leq 0$ , then set  $\theta_{t+1} = \theta^*$ .
  - b) If  $\Delta > 0$ , then set  $\theta_{t+1} = \theta^*$  with probability  $e^{-\Delta/\tau_j}$ , and  $\theta_{t+1} = \theta_t$  otherwise. This can be done as follows. Sample  $U \sim \text{Unif}(0, 1)$ . Then

$$\theta_{t+1} = \begin{cases} \theta^* & \text{if } U \leq e^{-\Delta/\tau_j} \\ \theta_t & \text{otherwise.} \end{cases}$$

3. Repeat steps 1 and 2 a total of  $m_j$  times.
4. Update  $\tau_{j+1} = \alpha(\tau_j)$ ,  $m_{j+1} = \beta(m_j)$  and move to stage  $j + 1$ , where  $\alpha$  and  $\beta$  are two deterministic functions that govern how to cool down the temperature and how long to stay at a given temperature.

Suppose  $I$  is an image consisting of “pixels”  $I(i, j)$ ,  $i, j = 1, \dots, N$ . The image is corrupted by noise  $Z(i, j)$  and only  $J = I + Z$  is observed. To reconstruct  $I$ , one way is to minimize a function

$$f(I) = \sum_{i,j} |J(i, j) - I(i, j)|^2 + K(I),$$

where  $K(I)$  is a function that has large values if  $I$  has many “irregularities”. The idea is that, as long as the noise is not too strong, the real image should be similar to  $J$ ; on the other hand, irregularities observed in  $J$  are likely to be due to noise and should be reduced. One way to use the simulated annealing to minimize  $f(I)$  is as follows. In each iteration, only one pixel of  $I$  can be updated. That means two  $I$  and  $I'$  are neighbors only when they are different at just one pixel. Then at each iteration, choose a pixel  $I(i, j)$  and select a candidate value for it. Update  $I(i, j)$  according to the rule in step 2 and move to the next iteration.

Some guidelines:

- R function

`optim`

can implement some simple versions of simulated annealing;

- the temperature  $\tau_j$  should slowly decrease to 0;
- the number  $m_j$  of iterations at each temperature  $\tau_j$  should be large and increasing in  $j$ ;
- reheating strategies that allow sporadic, systematic, or interactive temperature increases to prevent getting stuck in a local minimum at low temperatures can be effective.

# Genetic algorithms

First, each  $\theta \in \Theta$  is a string of alphabets:

$$\theta = (\theta_1, \dots, \theta_C),$$

where each  $\theta_i$  is a symbol from a finite set, such as  $\{0, 1\}$ ,  $\{0, 1, 2\}$ ,  $\{\text{'a'}, \text{'b'}, \dots\}$ .

Genetic algorithms regard the maximization of  $f(\theta)$  over  $\Theta$  as a process of natural selection, with  $f(\theta)$  being a measure of fitness and  $\theta$  the genetic code, or “chromosome”. It assumes that fitness is a result of some “good” pieces of  $\theta$  and by inheriting these pieces plus some mutations fitness is enhanced.

In a genetic algorithm, at each iteration  $t$ , at least two candidate solutions  $\theta_{t,1}, \dots, \theta_{t,P}$  have to be tracked. Each iteration consists of several steps.



- 1. Selection.** Randomly select from  $\theta_{t,1}, \dots, \theta_{t,P}$  to form a set of pairs. The selection should be based on a *fitness function*  $\phi(\theta)$ . In general, larger values of  $f(\theta)$  result in larger values of  $\phi(\theta)$ , and higher chance for  $\theta$  to be selected.

There are many selection mechanisms:

- a) select one parent  $\theta$  with probability proportional to  $\phi(\theta)$  and the other parent completely at random;
- b) select each parent independently with probability proportional to  $\phi(\theta)$ ;

$\phi$  must be selected carefully: using  $\phi(\theta_{t,i}) = f(\theta_{t,i})$  may result in rapid convergence into a local maximum. A common choice is

$$\phi(\theta_{t,i}) = \frac{2r_i}{P(P+1)}$$

where  $r_i$  is the *rank* of  $f(\theta_{t,i})$  in  $f(\theta_{t,1}), \dots, f(\theta_{t,P})$ .

- 2. Breeding.** For each pair,  $(\theta_a, \theta_b)$ , generate one or more “offspring”  $\theta' = c(\theta_a, \theta_b) \in \Theta$ , where  $c$  is a random operator. Typically,  $c$  is a “crossover”. If

$$\begin{aligned}\theta_a &= (\theta_{a1}, \dots, \theta_{aC}), \\ \theta_b &= (\theta_{b1}, \dots, \theta_{bC}),\end{aligned}$$

then a crossover works as follows,

- a) randomly choose a position  $1 \leq d \leq C$ ; and
- b) combine  $(\theta_{a1}, \dots, \theta_{ad})$  and  $(\theta_{b,d+1}, \dots, \theta_{bC})$  to form

$$\theta' = (\theta_{a1}, \dots, \theta_{ad}, \theta_{b,d+1}, \dots, \theta_{bC}).$$

If necessary, also take

$$\theta'' = (\theta_{a,d+1}, \dots, \theta_{aC}, \theta_{b1}, \dots, \theta_{bd})$$

to be another offspring.

- 3. Mutation.** Make some random modifications to each offspring. Typically, if an offspring chromosome is

$$\theta = (\theta_1, \dots, \theta_C)$$

then for  $i = 1, \dots, C$ , with a small probability  $0 < p < 1$  (mutation rate), change  $\theta_i$  to a different value.

The offspring produced at iteration  $t$  are taken as the candidate solutions for iteration  $t + 1$ .

# Initialization.

Usually the first generation consists of completely random individuals.

- Large values of the size of a generation,  $P$ , are preferred. For binary encoding of  $\theta$ , one suggestion is to have  $C \leq P \leq 2C$ , where  $C$  is the chromosome length. In most real applications, population sizes have ranged between 10 and 200.
- Mutation rates are typically very low, in the neighborhood of 1%.

# Termination.

A genetic algorithm is usually terminated after a maximum number of iterations. Alternatively, the algorithm can be stopped once the genetic diversity in the current generation is sufficiently low.

In many problems, like the traveling salesman problem, it is natural to write  $\theta$  as a permutation of  $1, 2, \dots, p$ .

- Standard crossover usually produces invalid sequences

For example, if

$$\theta_a = \{7, 5, 2, 6, 3, 1, 9, 4, 8\}$$

$$\theta_b = \{9, 1, 2, 3, 8, 6, 7, 5, 4\}$$

and if the crossover point is between 2nd and 3rd positions, then it produces

$$\{7, 5, 2, 3, 8, 6, 7, 5, 4\}$$

an invalid traveling route.

A remedy is order crossover: pick a number of positions from one parent and get the values at those positions, say  $i_1, \dots, i_s$ . Next identify those values from the 2nd parent. Suppose the set of values are found at  $j_1 < j_2 < \dots < j_s$ . Put  $i_1$  at  $j_1$ ,  $i_2$  at  $j_2$ ,  $\dots$ ,  $i_s$  at  $j_s$  while keeping the values of the 2nd on other locations unchanged.

Example: Consider parents (752631948) and (912386754) and random positions (4, 6, 7). The offsprings are (612389754) and (352671948).

The drawback of the operation is that it destroys some important structures of the parent routes, in particular, the links.

Edge-recombination crossover uses a different idea. It generates an offspring whose edges belong to those of the parent routes. By edge it means a link into or out of a city in a route. For example, the above two parents have the following links, the order of numbers in each link is unimportant.

$$(1, 2), (1, 3), (1, 9), (2, 3), (2, 5), (2, 6), (3, 6), (3, 8), \\ (4, 5), (4, 8), (4, 9), (5, 7), (6, 7), (6, 8), (7, 8)$$

First, choose one of the initial cities of the parents as the initial city of the offspring. At  $k$ -th step, if  $i_1, \dots, i_k$  have been chosen as the first  $k$  cities on the route, then choose among cities that are linked to  $i_k$  and are different from  $i_1, \dots, i_{k-1}$  as the  $(k+1)$ st city of the offspring.



## EM and MM Algorithms

# EM Optimizations

## What it is for

- inference when there is missing information
  - ▶ incomplete observations
  - ▶ auxiliary parameters that are unobserved but can facilitate inference on the main parameter.
  - ▶ Some important applications of the EM algorithm are in problems where one has what we might call pseudo missing data. We never had any chance of obtaining such data, but we could pretend they are missing and use EM algorithm to facilitate the computation of MLEs.

## Why it is important

- widely useful
- simple to implement
- numerically stable

Suppose a population consists of several sub-populations, each following a distribution  $p(x | \theta_k)$  and having population fraction  $q_k$ . Typically,

$$\mathbf{q} = (q_1, \dots, q_K), \quad \boldsymbol{\theta} = (\theta_1, \dots, \theta_K)$$

are unknown, and the goal is to estimate them.

Let  $x_1, \dots, x_n$  be a sample. If for each  $x_i$  we know it comes from the  $z_i^{\text{th}}$  sub-population, then, letting

$$\mathbf{x} = (x_1, \dots, x_n), \quad \mathbf{z} = (z_1, \dots, z_n),$$

the likelihood function of  $(\mathbf{q}, \boldsymbol{\theta})$  is

$$L(\mathbf{q}, \boldsymbol{\theta} | \mathbf{x}, \mathbf{z}) = \prod_{i=1}^n [q_{z_i} \times p(x_i | \theta_{z_i})].$$

However, in many cases,  $z$  is unobservable:

- cluster analysis
- the “sub-populations” may be constructs that facilitate inference or decision making, so they are nonexistent in reality.

The likelihood function then becomes

$$L(\mathbf{q}, \boldsymbol{\theta} | \mathbf{x}) = \sum_z L(\mathbf{q}, \boldsymbol{\theta} | \mathbf{x}, z).$$

The question is how to get the MLE of  $\mathbf{q}$  and  $\boldsymbol{\theta}$  efficiently.

In this example,  $\mathbf{y} = (\mathbf{x}, z)$  is the complete data. When only  $\mathbf{x}$  is observable,  $z$  is called the missing data.

Suppose a population follows a distribution  $p(y | \theta)$ , where  $\theta$  is the unknown parameter. Suppose a random sample is collected from  $p(y | \theta)$ . However, when an observation  $y_i$  is greater than a cut-off  $C$ , it is truncated and recorded as  $C$ , such as in clinical trials. In this case, we only know  $y_i \geq C$  but not its exact value. On the other hand, if  $y_i \leq C$ , then it is fully observed. Therefore, the observed data is

$$x_1 = \min(y_1, C), \dots, x_n = \min(y_n, C).$$

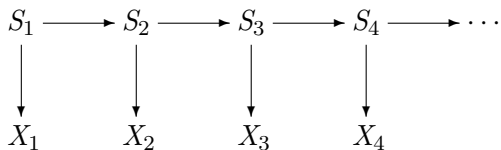
In this case,  $\mathbf{y} = (y_1, \dots, y_n)$  is the complete data and  $\mathbf{x} = (x_1, \dots, x_n)$ . The question is how to estimate  $\theta$  based on  $\mathbf{x}$ .

Suppose  $Z_1, \dots, Z_n$  and  $X_1, \dots, X_n$  are time series related by

$$X_k = f(Z_k | \theta),$$

where the transformation  $f$  is parametrized by  $\theta$ .  $Z_k$  may include not only variables of interest, but also some “noise” that interact with the variables. Suppose the joint distribution of  $Z_1, \dots, Z_n$  has the parametric form  $h(z_1, \dots, z_n | \gamma)$ , however, the values of  $\theta$  and  $\gamma$  are both unknown. If  $X_1, \dots, X_n$  are observed but  $Z_1, \dots, Z_n$  are only partially observed or completely hidden, how to estimate  $\theta$  and  $\gamma$ ?

The Hidden Markov Model (HMM) is a typical case. Let  $S_1, \dots, S_n$  be a discrete Markov chain with a finite number of states. Conditional on  $S_1, \dots, S_n$ , the observations  $X_1, \dots, X_n$  are independent, such that, given  $S_k = s$ ,  $X_k \sim N(\mu_s, \sigma_s^2)$ .



In this case,

$$Z_k = (S_k, W_k),$$

where  $W_1, \dots, W_n \sim N(0, 1)$  are independent of  $Z_1, \dots, Z_n$ , such that  $X_k = \mu_{S_k} + \sigma_{S_k} W_k$ . Typically,  $\mu_k$ ,  $\sigma_k$  and the transition probabilities of  $S_k$  are unknown and only  $X_k$  are observed.

Positron emission tomography (PET) is a tool to study the metabolic activity in organs. To generate a PET scan, some radioactive material is administered into the organ. Then the emissions of the material are recorded using a PET scanner, which consists of an array of detectors surrounding the patient's body. The region of interest is divided into "voxels". The number of photons  $Y_{ij}$  coming from voxel  $i$  and received by detector  $j$  is assumed to follow a Poisson distribution

$$Y_{ij} \sim \text{Poisson}(\theta_i a_{ij})$$

where the coefficient  $a_{ij}$  can be determined accurately beforehand, and  $\theta_i$  is the emission density of the radioactive material in voxel  $i$ , which is used to quantify the metabolic activity therein.



If for each pair of voxel and detector, we can observe  $Y_{ij}$ , then  $\theta_j$ 's can be estimated by MLE. However, in reality, each detector  $j$  receives photons from all the voxels, and hence only the sum

$$X_j = \sum_i Y_{ij}$$

is observable. The goal is to use  $\mathbf{X} = (X_j)$  instead of  $\mathbf{Y} = (Y_{ij})$  to estimate  $\boldsymbol{\theta} = (\theta_j)$ .

Suppose the complete data

$$\mathbf{Y} \sim p_{\mathbf{Y}}(\mathbf{y} | \boldsymbol{\theta})$$

while the observed data is  $\mathbf{X} = M(\mathbf{Y})$ , where  $M$  is a many-to-one mapping. Then

$$\underbrace{p_{\mathbf{X}}(\mathbf{x} | \boldsymbol{\theta})}_{L(\boldsymbol{\theta} | \mathbf{x})} = \int_{M(\mathbf{y})=\mathbf{x}} \underbrace{p_{\mathbf{Y}}(\mathbf{y} | \boldsymbol{\theta})}_{L(\boldsymbol{\theta} | \mathbf{y})} d\mathbf{y}.$$

If only  $\mathbf{X} = \mathbf{x}$  is observed, then the MLE  $\hat{\boldsymbol{\theta}}$  satisfies

$$L'(\hat{\boldsymbol{\theta}} | \mathbf{x}) = 0$$

i.e.,

$$\int_{M(\mathbf{y})=\mathbf{x}} L'(\hat{\boldsymbol{\theta}} | \mathbf{y}) d\mathbf{y} = 0.$$

Write log-likelihood function  $l(\boldsymbol{\theta} | \mathbf{y}) = \ln L(\boldsymbol{\theta} | \mathbf{y})$ , so that

$$[L(\boldsymbol{\theta} | \mathbf{y})]' = [e^{l(\boldsymbol{\theta} | \mathbf{y})}]' = [l(\boldsymbol{\theta} | \mathbf{y})]' e^{l(\boldsymbol{\theta} | \mathbf{y})} = [l(\boldsymbol{\theta} | \mathbf{y})]' L(\boldsymbol{\theta} | \mathbf{y}).$$

Therefore,

$$\int_{M(\mathbf{y})=x} L'(\hat{\boldsymbol{\theta}} | \mathbf{y}) d\mathbf{y} = 0$$

is the same as

$$\int_{M(\mathbf{y})=x} [l(\hat{\boldsymbol{\theta}} | \mathbf{y})]' L(\hat{\boldsymbol{\theta}} | \mathbf{y}) d\mathbf{y} = 0.$$

In principle, Newton's method can be used to find a solution. However, since  $\hat{\boldsymbol{\theta}}$  appears in two places in the integration, the implementation is often messy and because of that, numerically unstable.

We may try to “decouple” the two  $\theta$ ’s in the integration to get simpler iterations. An iterative procedure could go as follows. At step  $t$ , if we have  $\theta_t$ , then find  $\theta_{t+1}$  to solve

$$\int_{M(\mathbf{y})=x} [l(\theta_{t+1} | \mathbf{y})]' L(\theta_t | \mathbf{y}) \, d\mathbf{y} = 0,$$

or, to optimize the function

$$\int_{M(\mathbf{y})=x} l(\theta | \mathbf{y}) L(\theta_t | \mathbf{y}) \, d\mathbf{y}.$$

This is a function in  $\theta$ . Since  $\theta_t$  is different at each step, the function is different. If  $\theta_t$  converge, then the limit is a solution to the original MLE.

The explicit numerical integration in the above function is inconvenient. To replace it, note that, since  $\theta_t$  is fixed, the above way to find  $\theta_{t+1}$  is equivalent to optimizing

$$\int_{M(y)=x} l(\theta | y) \frac{L(\theta_t | y)}{L(\theta_t | x)} dy.$$

For any  $\theta_t$ ,

$$\frac{L(\theta_t | y)}{L(\theta_t | x)} = \frac{p_Y(y | \theta_t)}{p_X(x | \theta_t)} = p_{Y|X}(y | x, \theta_t),$$

the conditional density of  $Y$  given  $X = x$  under  $p_Y(y | \theta_t)$ . Then the integral is simply  $\mathbb{E}[l(\theta | Y) | x, \theta_t]$ . It turns out that in order to find suitable  $\theta_{t+1}$ , we need to maximize this function.

# EM (Expectation-Maximization) algorithm.

Define

$$Q(\boldsymbol{\theta} \mid \boldsymbol{\theta}') = \mathbb{E} [l(\boldsymbol{\theta} \mid \mathbf{Y}) \mid \mathbf{x}, \boldsymbol{\theta}'] = \mathbb{E} [\ln L(\boldsymbol{\theta} \mid \mathbf{Y}) \mid \mathbf{x}, \boldsymbol{\theta}'] .$$

- Initialize  $\boldsymbol{\theta}_0$ .
- At iteration  $t = 1, 2, \dots$ , with  $\boldsymbol{\theta}_t$  given, set  $\boldsymbol{\theta}_{t+1}$  to maximize  $Q(\boldsymbol{\theta} \mid \boldsymbol{\theta}_t)$ .

The “E step” of the algorithm refers to the computation of  $Q(\boldsymbol{\theta} \mid \boldsymbol{\theta}_t)$ , and the “M step” refers to the maximization of  $Q(\boldsymbol{\theta} \mid \boldsymbol{\theta}_t)$ . Stopping criteria are usually based upon  $|\boldsymbol{\theta}_{t+1} - \boldsymbol{\theta}_t|$  or  $|Q(\boldsymbol{\theta}_{t+1} \mid \boldsymbol{\theta}_t) - Q(\boldsymbol{\theta}_t \mid \boldsymbol{\theta}_t)|$ .

If  $\mathbf{Y} = (\mathbf{X}, \mathbf{Z})$ , where  $\mathbf{Z}$  is the missing data, then, given  $\mathbf{X} = \mathbf{x}$ , the only unknown part of  $\mathbf{Y}$  is  $\mathbf{Z}$ , so

$$\begin{aligned} Q(\boldsymbol{\theta} \mid \boldsymbol{\theta}') &= \mathbb{E} [\ln L(\boldsymbol{\theta} \mid \mathbf{Y}) \mid \mathbf{x}, \boldsymbol{\theta}'] \\ &= \mathbb{E} [\ln L(\boldsymbol{\theta} \mid \mathbf{x}, \mathbf{Z}) \mid \mathbf{x}, \boldsymbol{\theta}'] \\ &= \begin{cases} \sum_z p(z \mid \mathbf{x}, \boldsymbol{\theta}') \ln p(\mathbf{x}, z \mid \boldsymbol{\theta}), & \text{if } \mathbf{Z} \text{ is discrete} \\ \int p(z \mid \mathbf{x}, \boldsymbol{\theta}') \ln p(\mathbf{x}, z \mid \boldsymbol{\theta}) \, dz, & \text{if } \mathbf{Z} \text{ is continuous} \end{cases} \end{aligned}$$

# Gaussian mixture clustering

Suppose a population is a mixture of Gaussian distributions  $N(\boldsymbol{\mu}_k, \boldsymbol{\Sigma}_k)$  with population fractions  $q_k$ ,  $k = 1, \dots, K$ . The goal is to estimate  $\boldsymbol{\mu}_k$ ,  $\boldsymbol{\Sigma}_k$ ,  $q_k$ .

Let  $\mathbf{x}_1, \dots, \mathbf{x}_n \in \mathbb{R}^d$  be an iid sample from the population. Let  $z_i$  be the index of the Gaussian distribution from which  $\mathbf{x}_i$  is sampled. Suppose that

$$\mathbf{z} = (z_1, \dots, z_n)$$

cannot be observed. The observed data thus is

$$\mathbf{x} = (\mathbf{x}_1, \dots, \mathbf{x}_n).$$



Suppose at iteration  $t$  of the EM algorithm, one has obtained

$$\boldsymbol{\theta}_t = (\boldsymbol{\mu}_{t1}, \dots, \boldsymbol{\mu}_{tK}, \boldsymbol{\Sigma}_{t1}, \dots, \boldsymbol{\Sigma}_{tK}, q_{t1}, \dots, q_{tK}).$$

Then, since  $(\mathbf{x}_i, z_i)$  are independent, and  $z_i$  are discrete, for any candidate

$$\boldsymbol{\theta} = (\boldsymbol{\mu}_1, \dots, \boldsymbol{\mu}_K, \boldsymbol{\Sigma}_1, \dots, \boldsymbol{\Sigma}_K, q_1, \dots, q_K),$$

one has

$$\begin{aligned} Q(\boldsymbol{\theta} \mid \boldsymbol{\theta}_t) &= \sum_{\mathbf{z}} p(\mathbf{z} \mid \mathbf{x}, \boldsymbol{\theta}_t) \ln p(\mathbf{x}, \mathbf{z} \mid \boldsymbol{\theta}) \\ &= \sum_{i=1}^n \sum_{k=1}^K \underbrace{p(z_i = k \mid \mathbf{x}_i, \boldsymbol{\theta}_t)}_{w_{ik}} \ln p(z_i = k, \mathbf{x}_i \mid \boldsymbol{\theta}). \end{aligned}$$

To maximize  $Q(\boldsymbol{\theta} | \boldsymbol{\theta}_t)$ , note that  $w_{ik}$  has nothing to do with  $\boldsymbol{\theta}$  so it should be computed in the first place. By Bayes rule,

$$w_{ik} = \frac{p(z_i = k, \mathbf{x}_i | \boldsymbol{\theta}_t)}{p(\mathbf{x}_i | \boldsymbol{\theta}_t)} = \frac{p(z_i = k, \mathbf{x}_i | \boldsymbol{\theta}_t)}{\sum_{s=1}^K p(z_i = s, \mathbf{x}_i | \boldsymbol{\theta}_t)},$$

with

$$\begin{aligned} p(z_i = k, \mathbf{x}_i | \boldsymbol{\theta}_t) &= p(z_i = k | \boldsymbol{\theta}_t) p(\mathbf{x}_i | z_i = k, \boldsymbol{\theta}_t) \\ &= q_{tk} \phi(\mathbf{x}_i | \boldsymbol{\mu}_{tk}, \boldsymbol{\Sigma}_{tk}), \end{aligned} \quad (1)$$

where  $\phi(\mathbf{x} | \boldsymbol{\mu}, \boldsymbol{\Sigma})$  is the density of  $N(\boldsymbol{\mu}, \boldsymbol{\Sigma})$  at  $\mathbf{x}$ . Once  $w_{ik}$ 's are computed,

$$Q(\boldsymbol{\theta} | \boldsymbol{\theta}_t) = \sum_{k=1}^K \sum_{i=1}^n w_{ik} \ln p(z_i = k, \mathbf{x}_i | \boldsymbol{\theta}),$$

which looks a lot simpler.

As in (1),  $p(z_i = k, \mathbf{x}_i | \boldsymbol{\theta}) = q_k n(\mathbf{x}_i | \boldsymbol{\mu}_k, \boldsymbol{\Sigma}_k)$ . Since

$$n(\mathbf{x}_i | \boldsymbol{\mu}_k, \boldsymbol{\Sigma}_k) = \frac{(2\pi)^{-d/2}}{\sqrt{|\boldsymbol{\Sigma}_k|}} \exp \left\{ -\frac{1}{2} (\mathbf{x}_i - \boldsymbol{\mu}_k)' \boldsymbol{\Sigma}_k^{-1} (\mathbf{x}_i - \boldsymbol{\mu}_k) \right\},$$

then

$$\begin{aligned} Q(\boldsymbol{\theta} | \boldsymbol{\theta}_t) &= \sum_{k=1}^K \sum_{i=1}^n w_{ik} \ln[(2\pi)^{-d/2} q_k] - \frac{1}{2} \sum_{k=1}^K \sum_{i=1}^n w_{ik} \ln |\boldsymbol{\Sigma}_k| \\ &\quad - \frac{1}{2} \sum_{k=1}^K \sum_{i=1}^n w_{ik} (\mathbf{x}_i - \boldsymbol{\mu}_k)' \boldsymbol{\Sigma}_k^{-1} (\mathbf{x}_i - \boldsymbol{\mu}_k) \\ &= I_1 - \frac{I_2}{2} - \frac{I_3}{2}. \end{aligned}$$

From last page,

$$I_1 = \sum_{k=1}^K \sum_{i=1}^n w_{ik} \ln[(2\pi)^{-d/2} q_k], \quad I_2 = \sum_{k=1}^K \sum_{i=1}^n w_{ik} \ln |\boldsymbol{\Sigma}_k|$$

$$I_3 = \sum_{k=1}^K \sum_{i=1}^n w_{ik} (\mathbf{x}_i - \boldsymbol{\mu}_k)' \boldsymbol{\Sigma}_k^{-1} (\mathbf{x}_i - \boldsymbol{\mu}_k)$$

First, only  $I_3$  contains  $\boldsymbol{\mu}_k$ , which is a quadratic form. To *minimize* it, observe that  $I_3$  is the sum of  $K$  quadratic forms, each involving a single  $\boldsymbol{\mu}_k$

$$I_{3k} = \sum_{i=1}^n w_{ik} (\mathbf{x}_i - \boldsymbol{\mu}_k)' \boldsymbol{\Sigma}_k^{-1} (\mathbf{x}_i - \boldsymbol{\mu}_k).$$

We only need to find  $\boldsymbol{\mu}_k$  to minimize each  $I_{3k}$ . From the property of sample mean,  $\boldsymbol{\mu}_k$  must be the mean of a weighted sample  $\mathbf{x}_1, \dots, \mathbf{x}_n$ , each  $\mathbf{x}_i$  having weight  $w_{ik}$ . So

$$\boldsymbol{\mu}_k = \frac{\sum_{i=1}^n w_{ik} \mathbf{x}_i}{\sum_{i=1}^n w_{ik}}.$$

Next, only  $I_2$  and  $I_3$  contain  $\boldsymbol{\Sigma}_k$ , and  $I_2 + I_3$  is the sum of  $K$  terms of the following form, each involving a single  $\boldsymbol{\Sigma}_k$ ,

$$S_k = \sum_{i=1}^n w_{ik} \ln |\boldsymbol{\Sigma}_k| + \sum_{i=1}^n w_{ik} (\mathbf{x}_i - \boldsymbol{\mu}_k)' \boldsymbol{\Sigma}_k^{-1} (\mathbf{x}_i - \boldsymbol{\mu}_k)$$

and so we only need to find  $\boldsymbol{\Sigma}_k$  to minimize each  $S_k$ . Now that  $\boldsymbol{\mu}_k$  is equal to the weighted mean of  $\mathbf{x}_i$ , to *minimize*  $S_k$ ,  $\boldsymbol{\Sigma}_k$  must be the sample variance of the weighted sample  $\mathbf{x}_1, \dots, \mathbf{x}_n$ . So

$$\boldsymbol{\Sigma}_k = \frac{\sum_{i=1}^n w_{ik} (\mathbf{x}_i - \boldsymbol{\mu}_k)(\mathbf{x}_i - \boldsymbol{\mu}_k)'}{\sum_{i=1}^n w_{ik}}.$$

Finally, only  $I_1$  contains  $q_k$ . Since

$$\begin{aligned}
 I_1 &= \sum_{k=1}^K \sum_{i=1}^n w_{ik} \ln[(2\pi)^{-d/2} q_k] \\
 &= \sum_{k=1}^K \sum_{i=1}^n w_{ik} \ln[(2\pi)^{-d/2}] + \sum_{k=1}^K \sum_{i=1}^n w_{ik} \ln q_k \\
 &= -(d/2) \ln(2\pi) \sum_{k=1}^K \sum_{i=1}^n w_{ik} + \sum_{k=1}^K \left( \sum_{i=1}^n w_{ik} \right) \ln q_k
 \end{aligned}$$

it suffices to minimize

$$\sum_{k=1}^K \underbrace{\left( \sum_{i=1}^n w_{ik} \right)}_{W_k} \ln q_k$$

Note  $q_1, \dots, q_K$  must satisfy  $q_1 + \dots + q_K = 1$ . Therefore, the maximization can be obtained by finding a solution for  $\mathcal{L}'(q_1, \dots, q_K) = 0$ , i.e.

$$\frac{\partial \mathcal{L}(q_1, \dots, q_K)}{\partial q_k} = 0$$

where

$$\mathcal{L}(q_1, \dots, q_K) = \sum_{k=1}^K W_k \ln q_k - \lambda \left\{ \sum_{k=1}^K q_k - 1 \right\}$$

with  $\lambda$  a Lagrange multiplier. Then

$$q_k = \frac{W_k}{\sum_{k=1}^K W_k}.$$

Since for each  $i$ ,  $\sum_{k=1}^K w_{ik} = 1$ ,

$$q_k = \frac{1}{n} \sum_{i=1}^n w_{ik}.$$

The above steps consist a basic EM Gaussian clustering algorithm. It can be summarized as follows.

- Given

$$\boldsymbol{\theta}_t = (\boldsymbol{\mu}_{t1}, \dots, \boldsymbol{\mu}_{tn}, \boldsymbol{\Sigma}_{t1}, \dots, \boldsymbol{\Sigma}_{tK}, q_{t1}, \dots, q_{tK}),$$

for  $k = 1, \dots, K$ , compute

$$w_{ik} = p(z_i = k \mid \mathbf{x}_i, \boldsymbol{\theta}_t), \quad i = 1, \dots, n$$

then set  $q_{t+1,k}$  equal to the mean of  $w_{1k}, \dots, w_{nk}$ ,  $\boldsymbol{\mu}_{t+1,k}$  and  $\boldsymbol{\Sigma}_{t+1,k}$  equal to the weighted average and weighted covariance of  $\mathbf{x}_1, \dots, \mathbf{x}_n$ , with weights  $w_{1k}, \dots, w_{nk}$ , respectively.



For high dimensional data, there are many variants of the algorithm which put constraints on  $\mathbf{\Sigma}_k$  to improve computational and estimation efficiency. The strongest constraint (and computationally the simplest) is

$$\mathbf{\Sigma}_1 = \cdots = \mathbf{\Sigma}_K = \sigma^2 \mathbf{I},$$

with  $\sigma^2$  being the only parameter. Others include

- ①  $\mathbf{\Sigma}_k = \sigma_k^2 \mathbf{I}$ ,  $k = 1, \dots, K$ , with  $\sigma_k$  being the parameters;
- ②  $\mathbf{\Sigma}_k$  are diagonal matrices that may be equal to or different from each other;
- ③  $\mathbf{\Sigma}_1 = \sigma_k^2 \mathbf{\Sigma}$ , with the parameters being  $\sigma_k^2$  and  $\mathbf{\Sigma}$ .

If Newton's method is used to directly maximize the log likelihood function  $l(\boldsymbol{\theta} | \mathbf{x}) = \ln L(\boldsymbol{\theta} | \mathbf{x})$ , then one would have to differentiate the function

$$l(\boldsymbol{\theta} | \mathbf{x}) = \sum_{i=1}^n l_i(\boldsymbol{\theta})$$

where, for every  $i = 1, \dots, n$ ,

$$\begin{aligned} l_i(\boldsymbol{\theta}) &= \ln p(\mathbf{x}_i | \boldsymbol{\theta}) = \ln \left\{ \sum_{z_i} p(\mathbf{x}_i, z_i | \boldsymbol{\theta}) \right\} \\ &= \ln \left\{ \sum_{k=1}^K p(\mathbf{x}_k | z_i = k, \boldsymbol{\theta}) p(z_i = k | \boldsymbol{\theta}) \right\} \\ &= \ln \left\{ \sum_{k=1}^K \frac{(2\pi)^{-\frac{d}{2}} q_k}{\sqrt{|\boldsymbol{\Sigma}_k|}} e^{-\frac{1}{2}(\mathbf{x}_i - \boldsymbol{\mu}_k)' \boldsymbol{\Sigma}_k^{-1} (\mathbf{x}_i - \boldsymbol{\mu}_k)} \right\}. \end{aligned}$$

The iterations apparently are quite involved.

# Computation

## The R package

`mclust`

provides a large number of clustering methods, including EM based methods. Use

```
install.packages("mclust")
```

to install it on your computer if it is not yet installed. To use, first execute

```
library("mclust")
```

which loads all the functions in the package into an R session.

## `mclustBIC`

computes the Bayes information criterion associated with different Gaussian mixture models, which are characterized by number of clusters and constraints on the covariance matrices.

## `mclustModel`

estimate the parameters for the best model determined via `mclustBIC`.

## `mclustModelNames`

lists the names of models used in the `mclust` package.

## `mclust2Dplot`

plots 2-D data given parameters of a Gaussian mixture model for the data. Some MATLAB routines, such as `EM-GM.m` can perform some of the functions of the R package.

# EM for Exponential Families

Many important parametric distributions belong to an exponential family, which has the form

$$p_Y(\mathbf{y} | \boldsymbol{\theta}) = c_1(\mathbf{y}) c_2(\boldsymbol{\theta}) \exp \left\{ \boldsymbol{\theta}^T \mathbf{s}(\mathbf{y}) \right\},$$

where  $\boldsymbol{\theta}$  is a vector of parameters and  $\mathbf{s}(\mathbf{y})$  a vector of sufficient statistics.

- ① The normal densities  $N_d(\boldsymbol{\theta}, \boldsymbol{\Sigma})$  with  $\boldsymbol{\Sigma}$  being fixed consist an exponential family, because

$$\begin{aligned} p_Y(\mathbf{y} | \boldsymbol{\theta}) &= \frac{e^{-\frac{1}{2}(\mathbf{y}-\boldsymbol{\theta})^T \boldsymbol{\Sigma}^{-1}(\mathbf{y}-\boldsymbol{\theta})}}{(2\pi)^{d/2} |\boldsymbol{\Sigma}|^{1/2}} \\ &= \underbrace{\frac{e^{-\frac{1}{2}\mathbf{y}^T \mathbf{y}}}{(2\pi)^{d/2}}}_{c_1(\mathbf{y})} \underbrace{\frac{e^{-\frac{1}{2}\boldsymbol{\theta}^T \boldsymbol{\theta}}}{|\boldsymbol{\Sigma}|^{1/2}}}_{c_2(\boldsymbol{\theta})} \exp\left\{\boldsymbol{\theta}^T \underbrace{(\boldsymbol{\Sigma}^{-1} \mathbf{y})}_{s(\mathbf{y})}\right\}. \end{aligned}$$

- ① Poisson distribution  $\text{Poisson}(\theta)$ , where  $\theta > 0$ , has distribution function

$$p(y | \theta) = \frac{e^{-\theta} y^\theta}{y!}, \quad y = 0, 1, \dots$$

If we define

$$c_1(y) = \frac{1}{y!} \quad c_2(\theta) = e^{-\theta} \quad s(y) = \ln y,$$

then  $p(y | \theta) = c_1(y) c_2(\theta) e^{\theta s(y)}$ . Therefore,  $\text{Poisson}(\theta)$  also form an exponential family.

- ① Gamma distribution  $\text{Gamma}(a, r)$ , where  $a > 0$ ,  $r > 0$ , has density

$$p(y | a, r) = \begin{cases} \frac{e^{-y/r} y^{a-1}}{r^a \Gamma(a)} & y > 0 \\ 0 & y \leq 0. \end{cases}$$

Indeed, if we define  $b = 1/r$ ,  $\theta = (a, b)$ , and

$$c_1(y) = \frac{1}{y} \quad c_2(\theta) = c_2(a, b) = \frac{b^a}{\Gamma(\theta_1)} \quad s(y) = (\ln y, -y)$$

Then

$$c_1(y) c_2(\theta) \exp(\theta^T s(y)) = \frac{1}{y} \frac{1}{b^a \Gamma(a)} e^{a \ln y - by} = p(y | a, r).$$



- 1 Beta distribution  $\text{Beta}(a, b)$ , where  $a > 0$ ,  $b > 0$ , has density

$$p(y | a, b) = \begin{cases} \frac{y^{a-1}(1-y)^{b-1}}{B(a, b)}, & 0 < y < 1 \\ 0 & \text{otherwise} \end{cases}$$

If we let  $\theta = (a, b)$ ,

$$c_1(y) = \frac{1}{y(1-y)} \quad c_2(\theta) = \frac{1}{B(a, b)} \quad s(y) = (\ln y, \ln(1-y))$$

then

$$c_1(y)c_2(\theta)\exp(\theta^T s(y)) = \frac{1}{y(1-y)} \frac{1}{B(a, b)} e^{a \ln y + b \ln(1-y)} = p(y | a, b).$$

Suppose  $\mathbf{x} = M(\mathbf{y})$  is observed and we want to estimate  $\boldsymbol{\theta}$ . The EM algorithm in this case is significantly simplified, with no explicit evaluation of  $Q(\boldsymbol{\theta} | \boldsymbol{\theta}')$ . It works as follows. Given  $\boldsymbol{\theta}_t$ , set  $\boldsymbol{\theta}_{t+1}$  such that

$$\mathbb{E}[s(\mathbf{Y}) | \boldsymbol{\theta}_{t+1}] = \mathbb{E}[s(\mathbf{Y}) | \mathbf{x}, \boldsymbol{\theta}_t]. \quad (2)$$

To get the result, we start with the evaluation

$$\begin{aligned} Q(\boldsymbol{\theta} | \boldsymbol{\theta}_t) &= \int [\ln p_Y(\mathbf{y} | \boldsymbol{\theta})] p(\mathbf{y} | \mathbf{x}, \boldsymbol{\theta}_t) d\mathbf{y} \\ &= \int [\ln c_1(\mathbf{y}) + \ln c_2(\boldsymbol{\theta}) + \boldsymbol{\theta}^T \mathbf{s}(\mathbf{y})] p(\mathbf{y} | \mathbf{x}, \boldsymbol{\theta}_t) d\mathbf{y} \\ &= \int [\ln c_1(\mathbf{y})] p(\mathbf{y} | \mathbf{x}, \boldsymbol{\theta}_t) d\mathbf{y} \\ &\quad + \int [\ln c_2(\boldsymbol{\theta})] p(\mathbf{y} | \mathbf{x}, \boldsymbol{\theta}_t) d\mathbf{y} + \int \boldsymbol{\theta}^T \mathbf{s}(\mathbf{y}) p(\mathbf{y} | \mathbf{x}, \boldsymbol{\theta}_t) d\mathbf{y} \\ &= k(\boldsymbol{\theta}_t) + \ln c_2(\boldsymbol{\theta}) + \boldsymbol{\theta}^T \mathbb{E}[s(\mathbf{Y}) | \mathbf{x}, \boldsymbol{\theta}_t]. \end{aligned}$$

To maximize  $Q(\boldsymbol{\theta} | \boldsymbol{\theta}_t)$  as a function of  $\boldsymbol{\theta}$ , it suffices to solve

$$(\ln c_2(\boldsymbol{\theta}))' + \mathbb{E}[s(\mathbf{Y}) | \mathbf{x}, \boldsymbol{\theta}_t] = 0,$$

or

$$\frac{c_2'(\boldsymbol{\theta})}{c_2(\boldsymbol{\theta})} + \mathbb{E}[s(\mathbf{Y}) | \mathbf{x}, \boldsymbol{\theta}_t] = 0.$$

The following equality is basic for exponential families. For *any*  $\boldsymbol{\theta}$ ,

$$\frac{c_2'(\boldsymbol{\theta})}{c_2(\boldsymbol{\theta})} = -\mathbb{E}[s(\mathbf{Y}) | \boldsymbol{\theta}] = -\int s(\mathbf{y})p_{\mathbf{Y}}(\mathbf{y} | \boldsymbol{\theta}) \, d\mathbf{y}. \quad (3)$$

Assuming the result is correct for now, it is seen that  $\boldsymbol{\theta}_{t+1}$  should be set such that  $\mathbb{E}[s(\mathbf{Y}) | \boldsymbol{\theta}_{t+1}] = \mathbb{E}[s(\mathbf{Y}) | \mathbf{x}, \boldsymbol{\theta}_t]$ , i.e., as in (2).

To get (3), notice  $\int p_Y(\mathbf{y} | \boldsymbol{\theta}) = 1$  for any  $\boldsymbol{\theta}$ , i.e.

$$c_2(\boldsymbol{\theta}) \int c_1(\mathbf{y}) \exp\{\boldsymbol{\theta}^T \mathbf{s}(\mathbf{y})\} d\mathbf{y} = 1.$$

Then

$$\ln c_2(\boldsymbol{\theta}) = -\ln \left\{ \int c_1(\mathbf{y}) \exp\{\boldsymbol{\theta}^T \mathbf{s}(\mathbf{y})\} d\mathbf{y} \right\}.$$

Differentiate both sides to get

$$\begin{aligned} \frac{c'_2(\boldsymbol{\theta})}{c_2(\boldsymbol{\theta})} &= -\frac{\int c_1(\mathbf{y}) \mathbf{s}(\mathbf{y}) \exp\{\boldsymbol{\theta}^T \mathbf{s}(\mathbf{y})\} d\mathbf{y}}{\int c_1(\mathbf{y}) \exp\{\boldsymbol{\theta}^T \mathbf{s}(\mathbf{y})\} d\mathbf{y}} \\ &= -c_2(\boldsymbol{\theta}) \int c_1(\mathbf{y}) \mathbf{s}(\mathbf{y}) \exp\{\boldsymbol{\theta}^T \mathbf{s}(\mathbf{y})\} d\mathbf{y} \\ &= -\int \mathbf{s}(\mathbf{y}) p_Y(\mathbf{y} | \boldsymbol{\theta}) d\mathbf{y}, \end{aligned}$$

as claimed.

# Convergence of EM

We first show that each step of the algorithm increases the observed-data log likelihood  $l(\boldsymbol{\theta} | \boldsymbol{x}) = \ln p_X(\boldsymbol{x} | \boldsymbol{\theta})$ , that is

$$l(\boldsymbol{\theta}_{t+1} | \boldsymbol{x}) \geq l(\boldsymbol{\theta}_t | \boldsymbol{x}).$$

Once this is done, one can say  $l(\boldsymbol{\theta}_t | \boldsymbol{x})$  always climbs up to a local maximum of  $l(\boldsymbol{\theta} | \boldsymbol{x})$ .

- In general, there is no guarantee that the limit is the (global) maximum.
- EM has slower convergence than the Newton's method, sometime substantially.
- Nevertheless, the ease of implementation and the stable ascent of EM are often very attractive.

# Convergence of EM

By  $X = M(Y)$ ,  $p_Y(y|\theta) = p(x, y|\theta) = p_X(x|\theta)p_{Y|X}(y|x, \theta)$ , so

$$p_X(x|\theta) = \frac{p_Y(y|\theta)}{p_{Y|X}(y|x, \theta)}.$$

Because  $x$  is fixed,  $p_{Y|X}(y|x, \theta)$  only depends on  $\theta$ . To simplify notation, write  $f(y|\theta) = p_{Y|X}(y|x, \theta)$ . Then

$$\ln p_X(x|\theta) = \ln p_Y(y|\theta) - \ln f(y|\theta).$$

## Convergence of EM

Integrate both sides with respect to the density  $f(\mathbf{y} | \boldsymbol{\theta}_t)$ . From

$$\begin{aligned} & \int [\ln p_Y(\mathbf{y} | \boldsymbol{\theta})] f(\mathbf{y} | \boldsymbol{\theta}_t) \, d\mathbf{y} \\ &= \int [\ln p_Y(\mathbf{y} | \boldsymbol{\theta})] p_{Y|X}(\mathbf{y} | \mathbf{x}, \boldsymbol{\theta}_t) \, d\mathbf{y} = Q(\boldsymbol{\theta} | \boldsymbol{\theta}_t), \end{aligned}$$

It follows

$$\ln p_X(\mathbf{x} | \boldsymbol{\theta}_t) = Q(\boldsymbol{\theta} | \boldsymbol{\theta}_t) - H(\boldsymbol{\theta} | \boldsymbol{\theta}_t), \quad (4)$$

where

$$H(\boldsymbol{\theta} | \boldsymbol{\theta}_t) = \int [\ln f(\mathbf{y} | \boldsymbol{\theta})] f(\mathbf{y} | \boldsymbol{\theta}_t) \, d\mathbf{y}.$$



# Convergence of EM

To show  $\ln p_{\mathbf{X}}(\mathbf{x} | \boldsymbol{\theta}_{t+1}) \geq \ln p_{\mathbf{X}}(\mathbf{x} | \boldsymbol{\theta}_t)$ , by (4), their difference equals

$$[Q(\boldsymbol{\theta}_{t+1} | \boldsymbol{\theta}_t) - Q(\boldsymbol{\theta}_t | \boldsymbol{\theta}_t)] - [H(\boldsymbol{\theta}_{t+1} | \boldsymbol{\theta}_t) - H(\boldsymbol{\theta}_t | \boldsymbol{\theta}_t)].$$

By the M-step,  $Q(\boldsymbol{\theta}_{t+1} | \boldsymbol{\theta}_t) \geq Q(\boldsymbol{\theta}_t | \boldsymbol{\theta}_t)$ . On the other hand, by the Jensen's inequality,

$$\begin{aligned} H(\boldsymbol{\theta}_{t+1} | \boldsymbol{\theta}_t) - H(\boldsymbol{\theta}_t | \boldsymbol{\theta}_t) &= \int \ln \left[ \frac{f(\mathbf{y} | \boldsymbol{\theta}_{t+1})}{f(\mathbf{y} | \boldsymbol{\theta}_t)} \right] f(\mathbf{y} | \boldsymbol{\theta}_t) \, d\mathbf{y} \\ &\leq \ln \left[ \int \frac{f(\mathbf{y} | \boldsymbol{\theta}_{t+1})}{f(\mathbf{y} | \boldsymbol{\theta}_t)} f(\mathbf{y} | \boldsymbol{\theta}_t) \, d\mathbf{y} \right] = \ln \left[ \int f(\mathbf{y} | \boldsymbol{\theta}_{t+1}) \, d\mathbf{y} \right] = 0, \end{aligned}$$

thus showing the likelihood increases at each iteration.

# Variants of EM

- E step is hard: Monte Carlo (MCEM)
- M step is hard: Instead of choosing  $\theta_{t+1}$  to maximize  $Q(\theta | \theta_t)$ , one simply chooses any  $\theta_{t+1}$  for which

$$Q(\theta_{t+1} | \theta_t) > Q(\theta_t | \theta_t),$$

then the resulting algorithm is called generalized EM, or GEM.

- ▶ In GEM,  $l(\theta_t | x)$  is still increasing.
- ▶ An example: ECM algorithm, which replaces the M-step of the EM algorithm with several computationally simpler conditional maximization (CM) steps.

# Variance Estimation

- Louis' method (1982, JRSSB)
- Supplemented EM (SEM) algorithm (Meng and Rubin, 1991, JASA)
- Oakes' method (1999, JRSSB)
- Bootstrapping
- Empirical information
- Numerical Differentiation

## Example: Hidden Markov Model

- Hidden state  $\mathbf{s} = (s_t)$ ,  $t = 1, \dots, T$ , follows a Markov chain with initial state distribution  $\pi_k$ ,  $k = 1, \dots, M$ , and transition probability matrix  $\{a_{k,l} = \Pr(s_{t+1} = l | s_t = k)\}$ . (Note this matrix is row stochastic.)
- Given the state  $s_t$ , the observation  $u_t$  is independent of other observations and states
- Given state  $k$ , the observation follows distribution  $b_k(u)$ . For example,  $b_k(u|\gamma_k)$  can be  $N(\mu_k, \Sigma_k)$ ; or discrete.
- The observed data is  $\mathbf{u} = (u_t)$ ,  $t = 1, \dots, T$ .

# Likelihood

- Parameters:  $\theta$  contains  $\pi_k$ 's,  $a_{kl}$ 's,  $\gamma_k$ 's
- Missing (latent, unobserved) states:  $s_t$ ,  $t = 1, \dots, T$ .
- Observed data:  $u_t$ ,  $t = 1, \dots, T$
- Complete data:  $(\mathbf{s}, \mathbf{u})$ , where  $\mathbf{s} = (s_1, \dots, s_T)$  and  $\mathbf{u} = (u_1, \dots, u_T)$ .
- Complete data likelihood

$$\pi_{s_1} \prod_{t=2}^T a_{s_{t-1}, s_t} \prod_{t=1}^T b_{s_t}(u_t | \gamma_{s_t})$$

- Complete data log-likelihood

$$\log p(\mathbf{s}, \mathbf{u} | \theta) = \log \pi_{s_1} + \sum_{t=2}^T \log a_{s_{t-1}, s_t} + \sum_{t=2}^T \log b_{s_t}(u_t | \gamma_{s_t})$$

## E-Step

Taking expectation of  $\log p(\mathbf{s}, \mathbf{u}|\theta')$  with respect to  $p(\mathbf{s}|\mathbf{u}, \theta)$ :

$$\begin{aligned}
 & E\left(\log p(\mathbf{s}, \mathbf{u}|\theta')|\mathbf{u}, \theta\right) \\
 &= \sum_{\mathbf{s}} p(\mathbf{s}|\mathbf{u}, \theta) \log p(\mathbf{s}, \mathbf{u}|\theta') \\
 &= \sum_{k=1}^M L_k(1) \log \pi'_k + \sum_{t=2}^T \sum_{k=1}^M \sum_{l=1}^M H_{k,l}(t) \log a'_{k,l} \\
 &\quad + \sum_{t=1}^T \sum_{k=1}^M L_k(t) \log b_{s_t}(u_t|\gamma'_k)
 \end{aligned}$$

where

$$L_k(t) = \Pr(s_t = k|\mathbf{u}) = \sum_{\mathbf{s}: s_t=k} p(\mathbf{s}|\mathbf{u})$$

and

$$H_{k,l}(t) = \Pr(s_t = k, s_{t+1} = l|\mathbf{u}) = \sum_{\mathbf{s}: s_t=k, s_{t+1}=l} \Pr(\mathbf{s}|\mathbf{u}).$$

# M-step

If  $b_k$ 's are multivariate normals, the M-step has closed-form solution:

$$\begin{aligned}\mu_k &= \frac{\sum_{t=1}^T L_k(t) u_t}{\sum_{t=1}^T L_k(t)} \\ \Sigma_k &= \frac{\sum_{t=1}^T L_k(t) (u_t - \mu_k)(u_t - \mu_k)^\top}{\sum_{t=1}^T L_k(t)} \\ a_{k,l} &= \frac{\sum_{t=1}^{T-1} H_{k,l}(t)}{\sum_{t=1}^{T-1} L_k(t)} \\ \pi_k &= \frac{\sum_{t=1}^T L_k(t)}{\sum_{k=1}^M \sum_{t=1}^T L_k(t)}\end{aligned}$$

# Forward and Backward Probability

- This is for efficient computation of  $L_k(t)$ 's and  $H_{k,l}(t)$ 's
- Computation of order  $M^2T$  and memory requirement of order  $MT$ .
- Forward probability

$$\alpha_k(t) = \Pr(u_1, \dots, u_t, s_t = k)$$

- Backward probability

$$\beta_k(t) = \Pr(u_{t+1}, \dots, u_T | s_t = k)$$



# Forward and Backward Algorithms

- Forward probability recursion

$$\alpha_k(1) = \pi_k b_k(u_1), \quad 1 \leq k \leq M,$$

$$\alpha_k(t) = b_k(u_t) \sum_{l=1}^M \alpha_l(t-1) a_{l,k}, \quad 1 < t \leq T, \quad 1 \leq k \leq M.$$

- Backward probability recursion

$$\beta_k(T) = 1, \quad 1 \leq k \leq M,$$

$$\beta_k(t) = \sum_{l=1}^M a_{k,l} b_l(u_{t+1}) \beta_l(t+1), \quad 1 \leq t < T, \quad 1 \leq k \leq M.$$

# Fast Algorithm for the E-step Quantities

$$L_k(t) = \Pr(s_t = k | \mathbf{u}) = \frac{\alpha_k(t)\beta_k(t)}{p(\mathbf{u})}$$

$$H_{k,l}(t) = \Pr(s_t = k, s_{k+1} = l | \mathbf{u}) = \frac{\alpha_k(t)a_{k,l}b_l(u_{t+1})\beta_l(t+1)}{p(\mathbf{u})}$$

$$p(\mathbf{u}) = \sum_{k=1}^M \alpha_k(t)\beta_k(t), \quad \forall t$$

## Majorization-Minimization Algorithm

# Majorization-Minimization Algorithm

- The MM algorithm is not an algorithm, but a prescription for constructing optimization algorithms.
- The EM algorithm is a special case.
- An MM algorithm operates by creating a surrogate function that majorizes (minorizes) the objective function. When the surrogate function is optimized, the objective function is driven downhill (uphill).
- Majorization-Minimization, Minorization-Maximization.

# Majorization-Minimization Algorithm

- It may generate an algorithm that avoids matrix inversion.
- It can linearize an optimization problem.
- It can conveniently deal with equality and inequality constraints.
- It can solve a non-differentiable problem by iteratively solving smooth problems.

# Majorization-Minimization Algorithm

- A function  $g(\theta \mid \theta^s)$  is said to majorize the function  $f(\theta)$  at  $\theta^s$  if

$$f(\theta^s) = g(\theta^s \mid \theta^s),$$

and

$$f(\theta) \leq g(\theta \mid \theta^s)$$

for all  $\theta$ .

- Let

$$\theta^{s+1} = \arg \min_{\theta} g(\theta \mid \theta^s).$$

- It follows that

$$f(\theta^{s+1}) \leq g(\theta^{s+1} \mid \theta^s) \leq g(\theta^s \mid \theta^s) = f(\theta^s).$$

The strict inequality holds unless  $g(\theta^{s+1} \mid \theta^s) = g(\theta^s \mid \theta^s)$  and  $f(\theta^{s+1}) = g(\theta^{s+1} \mid \theta^s)$ .

- Therefore, by alternating between the majorization and the minimization steps, the objective function is monotonically decreasing and thus its convergence is guaranteed.

# Majorization-Minimization Algorithm

- We can use any inequalities to construct the desired majorized/minorized version of the objective function. There are several typical choices:

- ▶ Jensen's inequality: for a convex function  $f(x)$ ,

$$f(E(X)) \leq E(f(X)).$$

- ▶ Convexity inequality: for any  $\lambda \in [0, 1]$ ,

$$f(\lambda x_1 + (1 - \lambda)x_2) \leq \lambda f(x_1) + (1 - \lambda)f(x_2)$$

- ▶ Cauchy-Schwarz inequality.
- ▶ Supporting hyperplanes.
- ▶ Arithmetic-Geometric Mean Inequality.

Example: Penalized AFT model with induced smoothing



# AFT model

- Let  $\{T_i, C_i, \mathbf{x}_i\}$ ,  $i = 1, \dots, n$  be  $n$  independent copies of  $\{T, C, \mathbf{X}\}$ , where  $T_i$  and  $C_i$  are log-transformed failure time and log transformed censoring time,  $\mathbf{x}_i$  is a  $p \times 1$  covariate vector, and given  $\mathbf{X}$ ,  $C$  and  $T$  are independent.
- The Accelerated failure time model has the form

$$T_i = \mathbf{x}_i^T \boldsymbol{\beta} + \epsilon_i, \quad i = 1, \dots, n.$$

- The observed data are  $(Y_i, \delta_i, \mathbf{x}_i)$ , where  $Y_i = \min(T_i, C_i)$ ,  $\delta_i = I(T_i < C_i)$ .

# AFT model

- For a case-cohort study, The weight-adjusted estimating equation with induced smoothing is

$$\tilde{U}_n(\beta) = \frac{1}{n} \sum_{i=1}^n \sum_{j=1}^n \Delta_i (\mathbf{x}_i - \mathbf{x}_j) \Phi \left( \frac{e_j(\beta) - e_i(\beta)}{r_{ij}} \right) = 0,$$

where  $e_j(\beta) = Y_j - \mathbf{x}_j^T \beta$  and

$$r_{ij}^2 = n^{-1} (\mathbf{x}_i - \mathbf{x}_j)^T (\mathbf{x}_i - \mathbf{x}_j) = n^{-1} \|\mathbf{x}_i - \mathbf{x}_j\|^2.$$

- Define  $H(x) = x\Phi(x) + \phi(x)$ , and

$$\tilde{L}_n(\beta) = \frac{1}{n} \sum_{i=1}^n \sum_{j=1}^n \Delta_i r_{ij} H \left( \frac{e_j(\beta) - e_i(\beta)}{r_{ij}} \right).$$

Then solving the estimating equation  $\tilde{U}_n(\beta) = 0$  is equivalent to minimizing the objective function  $\tilde{L}_n(\beta)$ .

- The objective function  $\tilde{L}_n(\beta)$  is convex in  $\beta$ .

# AFT model

- We propose the penalized AFT method, which amounts to minimize

$$\tilde{P}L_n(\beta) = \frac{1}{n} \sum_{i=1}^n \sum_{j=1}^n \Delta_i r_{ij} H \left( \frac{e_j(\beta) - e_i(\beta)}{r_{ij}} \right) + \lambda P(\beta),$$

where  $P(\beta)$  is some penalty function and  $\lambda$  is the tuning parameter.

# AFT model

- To solve this problem, the key is to identify a surrogate function of  $\tilde{L}_n(\beta)$  that is easy to work with. In view of the form of the objective function, an immediate idea is to majorize  $H(x)$ .
- Because  $H(x)$  is twice differentiable and  $H''(x) = \phi(x) < \phi(0) \approx 0.4$  for any  $x \in \mathbb{R}$ , it follows that for any  $x^0 \in R$ ,

$$H(x | x^0) \leq G(x | x^0) \equiv H(x^0) + (x - x^0)\Phi(x^0) + \frac{\phi(0)}{2}(x - x^0)^2.$$

- Note that in a standard optimization method based on quadratic approximation, e.g., Newton-Raphson, the  $\phi(0)$  term is replaced with  $\phi(x^0)$ , which, however, does not guarantee the majorization.

# AFT model

- Write

$$\tilde{L}_n(\beta) = \frac{1}{n} \sum_{i=1}^n \sum_{j=1}^n w_{ij} H(a_{ij} + \mathbf{z}_{ij}^T \beta) + \lambda P(\beta),$$

where  $w_{ij} = \Delta_i r_{ij}$ ,  $a_{ij} = (Y_j - Y_i)/r_{ij}$  and  $\mathbf{z}_{ij} = (\mathbf{x}_i - \mathbf{x}_j)/r_{ij}$ . Then

$$\begin{aligned} \tilde{L}_n(\beta \mid \beta^s) &\leq \frac{1}{n} \sum_{i=1}^n \sum_{j=1}^n w_{ij} G(a_{ij} + \mathbf{z}_{ij}^T \beta \mid a_{ij} + \mathbf{z}_{ij}^T \beta^s) + \lambda P(\beta) \\ &= \frac{1}{2} \beta^T \mathbf{A} \beta - \mathbf{b}^{s^T} \beta + \lambda P(\beta) + \text{const} \\ &\equiv \tilde{L}_n^*(\beta \mid \beta^s) + \text{const}, \end{aligned}$$

where

$$\mathbf{A} = \frac{\phi(0)}{n} \sum_{i=1}^n \sum_{j=1}^n w_{ij} \mathbf{z}_{ij} \mathbf{z}_{ij}^T, \quad \mathbf{b}^s = \mathbf{A} \beta^s - \frac{1}{n} \sum_{i=1}^n \sum_{j=1}^n w_{ij} \Phi(a_{ij} + \mathbf{z}_{ij}^T \beta^s) \mathbf{z}_{ij}$$

- Minimizing  $\tilde{L}_n^*(\beta \mid \beta^s)$  is a Lasso-type problem, for which a coordinate descent algorithm can be applied.

# MM for Penalized AFT

Initialize  $\beta^0 \in \mathbb{R}^p$ . Compute

$$\mathbf{A} = \frac{\phi(0)}{n} \sum_{i=1}^n \sum_{j=1}^n w_{ij} \mathbf{z}_{ij} \mathbf{z}_{ij}^T.$$

Set  $k \leftarrow 0$ .

**repeat**

(1). Majorization step:

$$\mathbf{b}^{k+1} \leftarrow \mathbf{A} \beta^k - \frac{1}{n} \sum_{i=1}^n \sum_{j=1}^n w_{ij} \Phi(a_{ij} + \mathbf{z}_{ij}^T \beta^k) \mathbf{z}_{ij}.$$

(2). Minimization step:

Initialize  $\tilde{\beta}^0 = \beta^k$ .

**while**  $\|\tilde{\beta}^s - \tilde{\beta}^{s+1}\| / \|\tilde{\beta}^s\| \geq \epsilon$  and  $s < s_{max}$  **do**

(i) For  $j = 1, \dots, p$ ,

$$\tilde{\beta}_j^s \leftarrow \frac{1}{a_{jj}} \mathcal{S}(b_j^{k+1} - \mathbf{a}_j^T \tilde{\beta}^s + a_{jj} \tilde{\beta}_j^s, \lambda).$$

(ii)  $\tilde{\beta}^{s+1} \leftarrow \tilde{\beta}^s$ .

(iii)  $s \leftarrow s + 1$ .

**end while**

$$\beta^{k+1} \leftarrow \tilde{\beta}^{s+1}.$$

Set  $k \leftarrow k + 1$ .

**until** convergence, i.e.,  $\|\beta^{k+1} - \beta^k\| / \|\beta^k\| \leq \epsilon$ .

## Simulation Basics

# Sampling Random Variables

The entire area of random sampling is built up the fact that we can use computing device to generate, deterministically, long sequences of numbers that by many accounts are very similar to long sequences of iid random variables following  $\text{Unif}(0, 1)$ .

In our discussion, we assume that we can generate ideal sequence  $U_1, U_2, \dots$  iid  $\sim \text{Unif}(0, 1)$ . The question is how to use them to simulate other random variables for vectors.



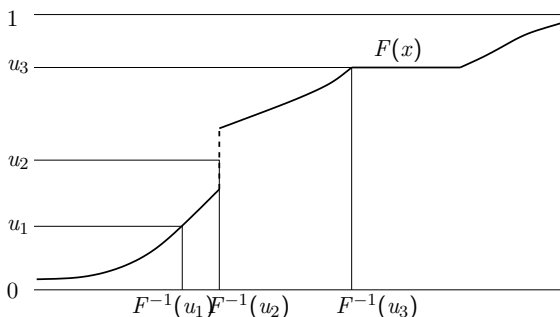
# Inverse transform method

This method only applies to random variables taking values in  $\mathbb{R}$ . Let  $F$  be a distribution function. Recall that  $F$  is nondecreasing and right-continuous with  $0 \leq F(x) \leq 1$ . However,  $F$  may be discontinuous at some  $x$  and may be constant in certain intervals.

For  $0 \leq u \leq 1$ , define

$$F^{-1}(u) = \inf\{x : F(x) \geq u\}.$$

The inverse can be thought as follows. Plot the graph of  $F$ . If  $F$  has a jump at  $x$ , the graph of  $F$  will have a gap at  $x$ . Fill the gap with a vertical line. Draw a horizontal line at height  $u$  starting from  $-\infty$ . Then  $F^{-1}(u)$  is the  $x$ -coordinate of the first point that the horizontal line meets with the gap-filled graph of  $F$ .



## Theorem

If  $U \sim \text{Unif}(0, 1)$ , then  $F^{-1}(U) \sim F$ .

## Proof.

For any  $x$ ,  $\mathbb{P}F^{-1}(U) \leq x = \mathbb{P}U \leq F(x) = F(x)$ . □

# Discrete Random Variable

Let  $X$  be a random variable taking values in  $\{c_1, \dots, c_n\}$ . Note that  $c_i$  may not be sorted or even not real numbers. Let

$$q_0 = 0, \quad q_i = \sum_{j=1}^i \mathbb{P}\{X = c_j\}, \quad i = 1, \dots, n$$

To sample  $X$ ,

- 1 sample  $U \sim \text{Unif}(0, 1)$ ;
- 2 find  $K \in \{1, \dots, n\}$ , such that  $q_{K-1} < U \leq q_K$ ;
- 3 set  $X = c_K$ .

□

# Exponential Variable

The exponential distribution with parameter  $\theta$ , usually denoted by  $\text{Exp}(\theta)$ , has

$$F(x) = 1 - e^{-x/\theta}, \quad x \geq 0.$$

Since  $F^{-1}(u) = -\theta \ln(1 - u)$  for  $u \in (0, 1)$ ,

$$-\theta \ln(1 - U) \sim \text{Exp}(\theta), \quad U \sim \text{Unif}(0, 1).$$

Since  $1 - U \sim \text{Unif}(0, 1)$ , there is also  $-\theta \ln U \sim \text{Exp}(\theta)$ .

A standard Brownian motion  $W(t)$  is a process with continuous sample paths that satisfies the following conditions

- (i)  $W(0) = 0$ ; and
- (ii) for  $0 \leq t_1 < \cdots < t_k$ , the increments

$$W(t_2) - W(t_1), \dots, W(t_k) - W(t_{k-1})$$

are independent and each  $W(t_i) - W(t_{i-1}) \sim N(0, t_i - t_{i-1})$ .

**Example** Let  $T$  be the time at which a standard Brownian motion  $W(t)$  attains its maximum over the time interval  $[0, 1]$ :

$$T = \inf\{t \in [0, 1] : W(t) = \sup_{x \in [0, 1]} W(x)\}.$$

Then  $T$  follows the *arcsine law* whose distribution is

$$F(x) = \frac{2}{\pi} \arcsin(\sqrt{x}), \quad x \in [0, 1]$$

and hence

$$\frac{1 - \cos(U\pi)}{2} \sim T.$$

□

## Proof.

For  $x \in [0, 1]$ ,

$$\begin{aligned}
 \mathbb{P}\{T \leq x\} &= \mathbb{P}\left\{\max_{t \in [0, x]} W(t) - W(x) \geq \max_{t \in [x, 1]} W(t) - W(x)\right\} \\
 &= \mathbb{P}\left\{\max_{t \in [-x, 0]} W(t) \geq \max_{t \in [0, 1-x]} W(t)\right\} \\
 &= \mathbb{P}\left\{\max_{t \in [0, x]} \tilde{W}(t) \geq \max_{t \in [0, 1-x]} W(t)\right\}
 \end{aligned}$$

where  $\tilde{W}$  is an independent copy of  $W$ . It is known (reflection principle) that  $\max_{t \in [0, x]} W(t) \sim x|Z|$ , where  $Z \sim N(0, 1)$ . It is also known that if  $Z$  and  $\tilde{Z}$  are iid  $N(0, 1)$ , then  $Z/\tilde{Z}$  is Cauchy with  $F(x) = \pi^{-1} \arctan x + 1/2$ . Therefore,  $\mathbb{P}\{T \leq x\} = \mathbb{P}\left\{(|Z/\tilde{Z}| \leq \sqrt{x/(1-x)})\right\} = \frac{2}{\pi} \arctan \sqrt{x/(1-x)} = \frac{2}{\pi} \arcsin(\sqrt{x})$ .  $\square$

**Example** Let

$$M = \sup_{t \in [0,1]} W(t)$$

Then, conditioning on  $W(1) = b$ ,  $M$  follows the *Rayleigh distribution*

$$\mathbb{P}\{M \leq x \mid W(1) = b\} = F_b(x) = 1 - e^{-2x(x-b)}, \quad x > \max(0, b).$$

Solving  $u = 1 - e^{-2x(x-b)}$  gives

$$x = \frac{b \pm \sqrt{b^2 - 2 \ln(1 - u)}}{2}$$

Since  $F_b(x) > 0$  only for  $x > b$ ,

$$F_b^{-1}(u) = \frac{b + \sqrt{b^2 - 2 \ln(1 - u)}}{2}.$$

Then

$$\frac{b + \sqrt{b^2 - 2 \ln U}}{2} \sim F_b.$$

□



### Proof.

Let  $T = \inf\{t : W(t) = x\}$ . For  $t \in (0, 1)$ ,  
 $f_{T|W(1)}(t | b) \propto f_{W(1)|T}(b | t)f_T(t)$ . By strong property of  $W$ ,

$$f_{W(1)|T}(b | t) = \frac{e^{-(b-x)^2/2(1-t)}}{\sqrt{2\pi(1-t)}}.$$

and by scaling  $\mathbb{P}\{T \leq t\} = \mathbb{P}\left\{\max_{s \in [0, t]} W(s) \geq x\right\} = \mathbb{P}\left\{M \geq x/\sqrt{t}\right\}$ .

□

## Sampling a conditional distribution

Let  $X \sim F$ . To sample  $X$  conditional on  $a < X \leq b$ , recall

$$\mathbb{P}\{a < X \leq b\} = F(b) - F(a).$$

In order for the sampling to make sense, assume  $F(a) < F(b)$ . Let  $A = F(a)$  and  $B = F(b)$ . Then

$$F^{-1}(A + (B - A)U)$$

follows the conditional distribution, since

$$\begin{aligned}\mathbb{P}\left\{F^{-1}(A + (B - A)U) \leq x\right\} &= \mathbb{P}\{A + (B - A)U \leq F(x)\} \\ &= \mathbb{P}\left\{U \leq \frac{F(x) - F(a)}{F(b) - F(a)}\right\} = \frac{F(x) - F(a)}{F(b) - F(a)}.\end{aligned}$$

The right hand side is just  $\mathbb{P}\{X \leq x \mid a < X \leq b\}$ .

**Example** Let  $c > 0$  and  $X \sim \text{Exp}(\theta)$ . To sample  $X$  conditional on  $X \geq c$ , recall the exponential distribution is memoryless in the sense that for any  $c > 0$ ,

$$X - c \sim \text{Exp}(\theta) \quad \text{given} \quad X > c.$$

Therefore, the conditional distribution can be sampled by  $-\theta \ln U + c$  with  $U \sim \text{Unif}(0, 1)$ . □

# Numerical evaluation of $F^{-1}$

If  $F^{-1}$  is not known explicitly, the inverse transform method is still applicable through numerical evaluation of  $F^{-1}$ . If  $F$  is continuous, computing  $F^{-1}(u)$  is equivalent to finding a root  $x$  of the equation  $F(x) - u = 0$ . For a distribution  $F$  with density  $f$ , Newton's method procedures a sequence

$$x_{n+1} = x_n - \frac{F(x_n) - u}{f(x_n)}$$

iteratively, given a starting point  $x_0$ . Under suitable conditions,  $x_n \rightarrow F^{-1}(u)$ .

R provides many sampling functions. In each of the following functions,  $n$  is the number of observations to be sampled.

- `runif(n, min=0, max=1)`:  $\text{Unif}(a, b)$ ;
- `rnorm(n, mean=0, sd=1)`:  $N(\mu, \sigma^2)$ ;
- `rpois(n, lambda)`:  $\text{Poisson}(\lambda)$

$$p_n = e^{-\lambda} \lambda^n / n!.$$

- `rbeta(n, shape1, shape2)`:  $\text{Beta}(\alpha, \beta)$

$$f(x) = x^{\alpha-1} (1-x)^{\beta-1} / B(\alpha, \beta), \quad 0 < x < 1;$$

note  $\alpha$  and  $\beta$  must be positive.

- `rgamma(n, shape, rate=1)`:  $\text{Gamma}(\alpha, 1/r)$

$$f(x) = r^\alpha x^{\alpha-1} e^{-rx} / \Gamma(\alpha). \quad x > 0;$$

note  $\alpha$  must be positive.

# Rejection sampling

The method was introduced by Von Neumann and is among the most widely applicable mechanisms. In standard statistical software, the method is implemented to randomly sample many important distributions, such as normal distributions, Gamma distributions, Beta distributions, Poisson distributions. Much more detail on specific developments can be found in the following books.

- 1 Wolfgang Hörmann, Josef Leydold, and Gerhard Derflinger, *Automatic nonuniform random variate generation*, Springer-Verlag, 2004.
- 2 Luc Devroye, *Nonuniform random variate generation*, Springer-Verlag, 1986.

Let  $\mathcal{X}$  be a discrete set or a Euclidean space. Suppose the goal is to sample from a target density  $f$  defined on  $\mathcal{X}$ . Suppose we know that  $f(x)$  is proportional to a function  $q(x)$ , i.e.,

$$f(x) = q(x)/C, \quad x \in \mathcal{X}.$$

In order for  $f$  to be a probability density,  $C$  must satisfy

$$C = \begin{cases} \int_{\mathcal{X}} q(x) \mathrm{d}x & \text{if } \mathcal{X} \text{ is continuous} \\ \sum_{x \in \mathcal{X}} q(x) & \text{if } \mathcal{X} \text{ is discrete} \end{cases}$$

Note that  $C$  is often intractable, especially when  $\mathcal{X}$  or  $q$  is complicated.

Let  $g$  be an “instrumental” density on  $\mathcal{X}$  from which we know how to generate samples (easily). Suppose  $g$  has the property that for some constant  $\alpha > 0$ ,

$$q(x) \leq \alpha g(x), \quad \text{for all } x \in \mathcal{X}.$$

The function  $\alpha g(x)$  is known as an *envelope*.

### Rejection sampling

- 1 Sample  $X \sim g$  and  $U \sim \text{Unif}(0, 1)$
- 2 If

$$U > \frac{q(X)}{\alpha g(X)}$$

then go to Step 1; otherwise return  $X$

The returned value is a random sample from the density  $f(x)$ .



Even when the complete form of  $f(x)$  is known, it is often still advantageous to use rejection sampling, by simply using  $q(x) = f(x)$ . As can be seen, the number of iterations until a sample is accepted follows a geometric distribution, and in each iteration, the probability of acceptance is

$$p_a = \frac{1}{\alpha} \int_{\mathcal{X}} q(x) \mathrm{d}x.$$

To reduce the expected number of iterations,  $p_a$  should be large. Equivalently,  $\alpha$  should be small. Because  $\alpha$  must satisfy  $q(x) \leq \alpha g(x)$ , given  $g(x)$ , the smallest possible  $\alpha$  is

$$\alpha = \sup \frac{q(x)}{g(x)}.$$

To further reduce  $\alpha$ , one has to choose appropriate  $g$ . Of course, one has to make sure  $g$  is easy to sample at the same time.

Sometimes, even though the instrument density  $g$  is easy to sample, it is tedious to write down the complete form of  $g(x)$ . Suppose  $g(x) = Ch(x)$ , where  $C$  is the known (but complicated) normalizing constant. Then the algorithm can be modified as follows. Let  $\beta$  be a constant such that

$$\frac{q(x)}{h(x)} \leq \beta$$

for all  $x$  with  $q(x) > 0$ . (Actually,  $\beta = \alpha/C$ .) Then in Step 2, replace  $q(X)/\alpha g(X)$  with

$$\frac{q(X)}{\beta h(X)}.$$

Sometimes it is easier to express  $f(x)$  as a function proportional to the sum of a set of functions,

$$f(x) \propto \sum_{k=1}^m q_k(x).$$

Superpose that we can find densities  $g_1(x), \dots, g_m(x)$  that are easy to sample, and constants  $\alpha_1, \dots, \alpha_m$ , such that

$$q_k(x) \leq \alpha_k g_k(x).$$

Then the rejection sampling can be modified as follows. Note that the complete form of each  $g_k(x)$  has to be known.

## Rejection sampling

- 1 Sample  $k$  from  $\{1, \dots, m\}$  with probabilities  $p_k \propto \alpha_k$ .
- 2 Sample  $X \sim g_k$  and  $U \sim \text{Unif}(0, 1)$
- 3 If

$$U > \frac{q_k(X)}{\alpha_k g_k(X)}$$

then go to Step 1; otherwise return  $X$

The returned value is a random sample from the density  $f(x)$ .

Proof: Let  $Y$  be a sampled value. To show that  $Y$  has density  $f$ , it suffices to show that for any set  $A \subset \mathcal{X}$ ,  $\mathbb{P}\{Y \in A\} = \int_A f(x) dx$ . From the procedure,

$$\begin{aligned}\mathbb{P}\{Y \in A\} &= \mathbb{P}\{X \in A \mid X \text{ is accepted}\} \\ &= \frac{\mathbb{P}\{X \in A \text{ and } X \text{ is accepted}\}}{\mathbb{P}\{X \text{ is accepted}\}}.\end{aligned}$$

Let  $C = \alpha_1 + \cdots + \alpha_m$ . Then

$$\begin{aligned}\mathbb{P}\{X \in A \text{ and } X \text{ is accepted}\} &= \sum_{k=1}^m p_k \mathbb{P}\{X \in A \text{ and } X \text{ is accepted} \mid X \text{ is from } g_k\} \\ &= C^{-1} \sum_{k=1}^m \alpha_k \mathbb{P}\{X \in A \text{ and } X \text{ is accepted} \mid X \text{ is from } g_k\}\end{aligned}$$

For each  $k$ , let

$$r(x) = \frac{q_k(x)}{\alpha_k g_k(x)}.$$

Then

$$\begin{aligned} & \mathbb{P}\{X \in A \text{ and } X \text{ is accepted} \mid X \text{ is from } g_k\} \\ &= \mathbb{P}\{X \in A, U \leq r(X) \mid X \text{ is from } g_k\} \end{aligned}$$

Given  $X = x$ , the probability that  $U \leq r(X)$  is simply  $r(x)$ . Therefore,

$$\begin{aligned} \mathbb{P}\{X \in A, U \leq r(X) \mid X \text{ is from } g_k\} &= \int \mathbf{1}\{x \in A\} r(x) g_k(x) \, dx \\ &= \int_A \frac{q_k(x)}{\alpha_k g_k(x)} g_k(x) \, dx \\ &= \frac{1}{\alpha_k} \int_A q_k(x) \, dx. \end{aligned}$$

Note the factor  $1/\alpha_k$ . As a result,

$$\mathbb{P}\{X \in A \text{ and } X \text{ is accepted}\} = C^{-1} \sum_{k=1}^m \int_A q_k(x) \, dx.$$

In particular, applying the formula to  $A = \mathcal{X}$ ,

$$\mathbb{P}\{X \text{ is accepted}\} = C^{-1} \sum_{k=1}^m \int q_k(x) \, dx,$$

As a result,

$$\mathbb{P}\{Y \in A\} = \frac{\sum_{k=1}^m \int_A q_k(x) \, dx}{\sum_{k=1}^m \int q_k(x) \, dx} = \int_A f(x) \, dx.$$

## Example: Beta Variable

For  $a, b > 0$ ,  $\text{Beta}(a, b)$  has a density

$$f(x) = \begin{cases} \frac{x^{a-1}(1-x)^{b-1}}{B(a, b)}, & 0 \leq x \leq 1 \\ 0 & \text{elsewhere} \end{cases}$$

and can be sampled by using the fact for independent  $Y \sim \text{Gamma}(a, 1)$  and  $Z \sim \text{Gamma}(b, 1)$ ,

$$\frac{Y}{Y + Z} \sim \text{Beta}(a, b).$$

One can try some easy alternatives without sampling Gamma distributions.



# Method 1

This method works only when  $a, b \geq 1$ . If  $a = b = 1$ , then  $\text{Beta}(a, b)$  is simply  $\text{Unif}(0, 1)$ . Let  $a > 1$  and  $b > 1$ . Then  $f(x)$  is maximized

$$x_0 = \frac{a-1}{a+b-2}.$$

Choose  $g$  to be the uniform density on  $[0, 1]$ . Then the optimal  $\alpha$  is

$$\alpha = \sup \frac{f(x)}{g(x)} = f(x_0).$$

The algorithm is implemented as follows.

- Draw  $X, U \sim \text{Unif}(0, 1)$  until  $f(x_0)U \leq f(X)$  and then return  $X$ .

## Method 2

Let  $U \sim \text{Unif}(0, 1)$ . Observe that for any  $a, b > 0$ ,

$$U^{1/a} \sim \text{Beta}(a, 1), \quad 1 - U^{1/b} \sim \text{Beta}(1, b).$$

The densities of the distributions are

$$f_a(x) = ax^{a-1}, \quad f_b(x) = b(1-x)^{b-1}.$$

Therefore, if  $g(x)$  is a mixture of  $f_a(x)$  and  $f_b(x)$ ,

$$g(x) = (1 - \lambda)f_a(x) + \lambda f_b(x),$$

where  $\lambda \in (0, 1)$ , then  $f(x)/g(x)$  is bounded.

To sample from  $g(x)$ , one can run the following steps.

- 1 Sample  $U$  and  $V \sim \text{Unif}(0, 1)$
- 2 If  $V > \lambda$ , then return  $X = U^{1/a}$ ; otherwise return  $X = 1 - U^{1/b}$ .  $\square$

# Conditional distributions

Let  $A$  be a subset of  $\mathcal{X}$ . To sample  $X$  conditional on  $X \in A$  one can use the following crude rejection sampling procedure

- Sample  $X$  until  $X \in A$  and then return  $X$ .

However, sometimes more carefully designed rejection sampling is needed to make efficient sampling.

## Example: Sampling from Normal Tails

Let  $c > 0$ . When  $c$  is large, to sample  $X$  from  $N(0, 1)$  conditional on  $X \geq c$ , the simple rejection sampling is very inefficient. To improve the efficiency, an exponential distribution can be used as the envelope. It suffices to sample  $X - c$ . Given  $X \geq c$ , the conditional density of  $X - c$  for  $X \sim N(0, 1)$  is

$$f(x) = \frac{\phi(x + c)}{1 - \Phi(c)}, \quad x \geq 0,$$

where  $\phi(x)$  is the standard normal density.

On the other hand, for  $X \sim \text{Exp}(\lambda)$ , the conditional density of  $X - c$  is

$$g(x) = \lambda e^{-\lambda x}.$$

Given  $\lambda$ , the optimal

$$\alpha = \sup_{x \geq 0} \frac{\phi(x + c)/[1 - \Phi(c)]}{\lambda e^{-\lambda x}} = \frac{\exp(\frac{1}{2}\lambda^2 - \lambda c)}{\sqrt{2\pi}\lambda[1 - \Phi(c)]}.$$

The value of  $\lambda$  that minimizes the last expression is

$$\lambda = \frac{c + \sqrt{c^2 + 4}}{2}.$$

Then the optimal

$$p_a = \frac{1}{\alpha} = \frac{\sqrt{\pi}(c + \sqrt{c^2 + 4})[1 - \Phi(c)]}{\sqrt{2}e}.$$

□

# Gamma distributions

For  $r, \lambda > 0$ ,  $\text{Gamma}(r, \lambda)$  has a density

$$f(x) = \frac{x^{r-1} e^{-x/\lambda}}{\lambda^r \Gamma(r)}, \quad x > 0.$$

- ❶  $\text{Exp}(\lambda) = \text{Gamma}(1, \lambda)$ . If  $U \sim \text{Unif}(0, 1)$ , then  $-\lambda \ln U \sim \text{Exp}(\lambda)$ .
- ❷  $\chi_k^2 = \text{Gamma}(\frac{k}{2}, \frac{1}{2})$ .
- ❸ If  $X \sim \text{Gamma}(r, \lambda)$  and  $c > 0$ , then  $cX \sim \text{Gamma}(r, c\lambda)$ .
- ❹ If  $X_i \sim \text{Gamma}(r_i, \lambda)$ ,  $i = 1, \dots, n$ , are independent, then

$$X_1 + \dots + X_n \sim \text{Gamma}(r_1 + \dots + r_n, \lambda).$$

Some simple ways to sample special Gamma distributions.

- ❶ If  $r$  is an integer, then  $\text{Gamma}(r, \lambda)$  can be sampled by

$$-\lambda \sum_{i=1}^r \ln U_i$$

where  $U_1, \dots, U_r$  are iid  $\sim \text{Unif}(0, 1)$

- ❷ If  $r$  is a half integer  $n/2$ , then  $\text{Gamma}(r, \lambda)$  can be sampled by

$$2\lambda \sum_{i=1}^n X_i^2$$

where  $X_1, \dots, X_n$  are iid  $\sim N(0, 1)$ .

Note that the above methods take longer time as  $r$  increases, because they require sampling of more and more variables.

There are many rejection sampling algorithms for Gamma distributions  $\text{Gamma}(r, 1)$  that are *universally fast*. That is, for such an algorithm, one can find a constant  $M > 0$ , such that for any  $r > 0$ , the expected number of iterations is no greater than  $M$ .

Typically, the cases  $r \geq 1$  and  $r < 1$  are treated separately. For  $r \geq 1$ , one universally fast algorithm, which is due to Best (1978), is as follows. Let

$$b = r - 1, \quad c = 3r - 3/4.$$

- ① Draw  $U, V \text{ iid } \sim \text{Unif}(0, 1)$  and set  $W = U(1 - U)$ ,  
 $Y = \sqrt{c/W}(U - 1/2)$ ,  $X = b + Y$ .
- ② Set  $Z = 64W^3V^2$ . If  $X \geq 0$ , and either  $Z \leq 1 - 2Y^2/X$  or  $\ln(Z) \leq 2[b \ln(X/b) - Y]$ , then return  $X$ ; otherwise go to Step 1.



For  $r \leq 1$ , the following algorithm of Ahrens and Dieter (1974) is universally fast. Let

$$p = \frac{e}{r + e}$$

and the instrument density

$$g(x) = \begin{cases} prx^{r-1}, & x \in [0, 1] \\ (1 - p)e^{-x+1} & x > 1. \end{cases}$$

That is,  $g$  is a mixture of  $\text{Beta}(r, 1)$  and  $\text{Exp}(1)$  conditional on  $(1, \infty)$ . Then keep sampling  $X \sim g$  and  $U \sim \text{Unif}(0, 1)$  until  $U \leq f(X)/[Cg(X)]$ , where  $f(X) = x^{r-1}e^{-x}/\Gamma(r)$  is the density of  $\text{Gamma}(r, 1)$  and  $C$  can be set equal to 1.39. In principle, the expected number of iterations is no greater than 1.39.

## Box–Müller method.

Let  $\Phi$  be the distribution function of  $N(0, 1)$ . To sample  $X \sim N(0, 1)$ , a straightforward method is to apply  $\Phi^{-1}(U)$  with  $U \sim \text{Unif}(0, 1)$ . There are two other often used methods.

This method does not need rejection sampling. However, the computation of sine and cosine functions can be time consuming.

- 1 Sample  $U_1, U_2$  iid  $\sim \text{Unif}(0, 1)$
- 2 Set  $R = \sqrt{-2 \ln U_1}$
- 3 Return

$$Z_1 = R \cos(2\pi U_2), \quad Z_2 = R \sin(2\pi U_2).$$

Then  $Z_1, Z_2$  are iid  $\sim N(0, 1)$ .

Proof: For any  $r > 0$ ,

$$\begin{aligned}\mathbb{P}\{R \leq r\} &= \mathbb{P}\{-2 \ln U_1 \leq r^2\} \\ &= \mathbb{P}\{U_1 \geq e^{-r^2/2}\} = 1 - e^{-r^2/2}\end{aligned}$$

So  $R$  has density  $re^{-r^2/2}$ . Recall that if  $\mathbf{X}$  is a random vector of  $d$  dimension and has a density  $f(\mathbf{x})$ , while

$$\phi(\mathbf{x}) = (\phi_1(\mathbf{x}), \dots, \phi_d(\mathbf{x}))$$

is a 1-to-1 differentiable function, then  $\mathbf{Y} = \phi(\mathbf{X})$  is a random vector with density

$$g(\mathbf{y}) = \frac{f(\mathbf{x})}{|\det[\phi'(\mathbf{x})]|}, \text{ with } \mathbf{x} = \phi^{-1}(\mathbf{y}).$$

Let  $\mathbf{X} = (R, U_2)$  and

$$\phi(\mathbf{x}) = (r \cos(2\pi u), r \sin(2\pi u)).$$

Then  $\mathbf{Y} = (Z_1, Z_2) = \phi(\mathbf{X})$ .

First,  $\mathbf{X} = (R, U_2)$  has a joint density  $f(r, u) = re^{-r^2/2}$ . Second,

$$\phi'(\mathbf{x}) = \begin{pmatrix} \cos(2\pi u) & -2\pi r \sin(2\pi u) \\ \sin(2\pi u) & 2\pi r \cos(2\pi u) \end{pmatrix}$$

and so  $|\det[\phi'(\mathbf{x})]| = 2\pi r$ . So  $\mathbf{Y} = (Z_1, Z_2)$  has density

$$g(z_1, z_2) = \frac{re^{-r^2/2}}{2\pi r} = \frac{e^{-r^2/2}}{2\pi},$$

with  $(r, u)$  satisfying

$$z_1 = r \cos(2\pi u_2), \quad z_2 = r \sin(2\pi u_2).$$

Then  $r^2 = z_1^2 + z_2^2$  and so

$$g(z_1, z_2) = \frac{e^{-(z_1^2 + z_2^2)/2}}{2\pi},$$

showing  $(Z_1, Z_2)$  are iid  $\sim N(0, 1)$ .

## Marsaglia–Bray method.

Let  $D$  be the disk centered at  $(0, 0)$  with radius 1. Then the method goes as follows:

- 1 use rejection sampling to sample  $(X, Y)$  from  $\text{Unif}(D)$
- 2 set  $R = \sqrt{X^2 + Y^2}$  and

$$T = \frac{2\sqrt{-\ln R}}{R}$$

- 3 return  $Z_1 = TX$ ,  $Z_2 = TY$

Again,  $Z_1$  and  $Z_2$  are iid  $\sim N(0, 1)$ .

# Multivariate normal distributions

Let  $\boldsymbol{\mu} \in \mathbb{R}^d$  and  $\boldsymbol{\Sigma} \in \mathbb{R}^{d \times d}$  be nonnegative definite. Recall that for  $\mathbf{X} \sim N(\boldsymbol{\mu}, \boldsymbol{\Sigma})$  and  $\mathbf{A} \in \mathbb{R}^{n \times d}$ ,

$$\mathbf{A}\mathbf{X} \sim N(\mathbf{A}\boldsymbol{\mu}, \mathbf{A}\boldsymbol{\Sigma}\mathbf{A}^t).$$

Based on this, to sample  $N(\boldsymbol{\mu}, \boldsymbol{\Sigma})$ , a method is to first compute the *Cholesky factorization* of  $\boldsymbol{\Sigma}$ , i.e.,

$$\boldsymbol{\Sigma} = \mathbf{A}\mathbf{A}^t,$$

with  $\mathbf{A}$  being lower triangular, and then sample

$$\boldsymbol{\mu} + \mathbf{A}\mathbf{Z},$$

with all the coordinates of  $\mathbf{Z}$  being iid  $\sim N(0, 1)$ .

## Conditioning formula

Let  $\mathbf{X}_1 \in \mathbb{R}^m$ ,  $\mathbf{X}_2 \in \mathbb{R}^n$  be jointly normal, such that

$$\mathbb{E}\mathbf{X}_i = \boldsymbol{\mu}_i, \quad \text{Cov}(\mathbf{X}_i, \mathbf{X}_j) = \boldsymbol{\Sigma}_{ij}.$$

If  $\boldsymbol{\Sigma}_{22}$  is of full rank, then, conditional on  $\mathbf{X}_2 = \mathbf{x}$ ,  $\mathbf{X}_1 \sim N(\boldsymbol{\mu}_c, \boldsymbol{\Sigma}_c)$ , where

$$\begin{aligned}\boldsymbol{\mu}_c &= \boldsymbol{\mu}_1 + \boldsymbol{\Sigma}_{12}\boldsymbol{\Sigma}_{22}^{-1}(\mathbf{x} - \boldsymbol{\mu}_2), \\ \boldsymbol{\Sigma}_c &= \boldsymbol{\Sigma}_{11} - \boldsymbol{\Sigma}_{12}\boldsymbol{\Sigma}_{22}^{-1}\boldsymbol{\Sigma}_{21}.\end{aligned}$$

Note that  $\boldsymbol{\Sigma}_c$  is independent of  $\mathbf{x}$ .

# Squeezed rejection sampling

Recall that in the ordinary rejection sampling of  $f(x)$ , we have  $q(x)$ ,  $g(x)$ , and  $\alpha$  available, such that  $f(x) \propto q(x)$ , where the proportionality factor may be intractable, and  $g(x)$  a density function such that  $q(x) \leq \alpha g(x)$ . Every time  $X \sim g$  is drawn,  $q(X)$  needs to be evaluated. If the evaluation is computationally costly, it slows down the simulation.

- 1 Squeezed rejection sampling can be used to preempt the evaluation of  $q(X)$ , hence improving simulation speed.
- 2 The sampling is still exact.

Choose a *squeezing function*  $s$  that is easy to compute, while satisfying

$$s(x) \leq q(x) \text{ for all } x$$



## Squeezed rejection sampling

- 1 Sample  $X \sim g$ ,  $U \sim \text{Unif}(0, 1)$ , and compute  $z = \alpha U g(X)$
- 2 If  $z \leq s(X)$ , then return  $X$   
otherwise, if  $z \leq q(X)$ , then return  $X$   
otherwise, go to Step 1.

In contrast, the ordinary rejection sampling directly goes from Step 1 to 3. Because of Step 2, the proportion of iterations in which evaluation of  $q$  is avoided is

$$\mathbb{P} \left\{ U \leq \frac{s(X)}{\alpha g(X)} \right\} = \frac{1}{\alpha} \int s(x) dx.$$

## Adaptive rejection sampling

The method developed by Gilks and Wild works well for density functions on an interval of  $\mathbb{R}$  which are continuous, differentiable, and *log-concave*. Log-concave densities are common in applications, for example, all  $N(\mu, \sigma^2)$ ,  $\text{Gamma}(\alpha, \beta)$  with  $\alpha \geq 1$ , and  $\text{Beta}(\alpha, \beta)$  with  $\alpha, \beta \geq 1$  are log-concave.

Suppose the density  $f$  is a density defined on  $(a, b)$ , where  $-\infty \leq a < b \leq \infty$ , such that

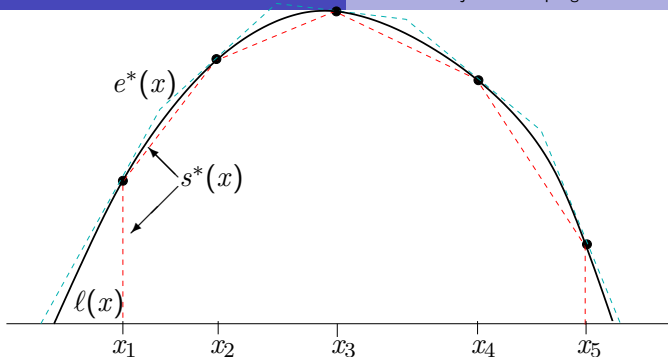
$$f(x) \propto q(x) = e^{\ell(x)}$$

where  $\ell(x)$  is concave. The idea of the method is to construct two piecewise linear functions  $e^*(x)$  and  $s^*(x)$ , such that

$$s^*(x) \leq \ell(x) \leq e^*(x)$$

Then  $e(x) = \exp(e^*(x))$  and  $s(x) = \exp(s^*(x))$  are envelope and squeezing functions respectively,

$$s(x) \leq q(x) \leq e(x).$$



The two piecewise linear functions  $s^*$  and  $e^*$  are shown above. First, select a set of points  $a < x_1 < \dots < x_n < b$ . Define  $s^*(x)$  to be  $-\infty$  outside  $[x_1, x_n]$  and to be linear on each  $[x_i, x_{i+1}]$ , such that its graph is the line segment connecting  $(x_i, \ell(x_i))$  and  $(x_{i+1}, \ell(x_{i+1}))$ . To construct  $e^*(x)$ , draw tangent lines to the graph of  $\ell$  at each  $(x_i, \ell(x_i))$  and cut off the parts that are beyond the intersection points of the tangent lines. Since  $\ell(x)$  is concave, it is guaranteed that  $s^*(x) \leq \ell(x) \leq e^*(x)$ . The “corner” points of  $e^*(x)$  can be derived in close form. We omit the detail for brevity.

Clearly,  $s(x) = \exp(s^*(x))$  can be used directly as a squeezing function. On the other hand, the  $e(x)$  is not a density function. Instead, letting

$$C = \int e(x) dx,$$

$g(x) = e(x)/C$  is a density function and satisfies  $q(x) \leq Cg(x)$ . The rest of the method is then the standard squeezed rejection sampling.

To see why the method often works so well, note the following.

- First,  $g(x)$  is an exponential function on each of the intervals  $(-\infty, x_1)$ ,  $[x_1, x_2]$ ,  $\dots$ ,  $[x_{n-1}, x_n]$ ,  $(x_n, \infty)$ . Therefore, it can be sampled using the sampling for exponential distributions, which is easy.
- Second,  $C$  can be evaluated easily as well, again because  $e(x)$  is piecewise exponential. Third, since  $s^*(x)$  is piecewise linear, it is very easy to evaluate. As a result,  $s(x) = \exp(s^*(x))$  is also easy to evaluate.

## Sampling importance resampling (SIR)

Like rejection sampling, SIR uses an auxiliary density  $g$ , this time called an importance sampling function. After getting a (large) sample from  $g$ , it is resampled according to appropriate weights. This *approximately* generates a sample from the target density  $f$ .

The weights are called the standardized importance weights. Given  $X_1, \dots, X_m \sim g$ , the weight for  $X_i$  is  $w(X_i) \propto f(X_i)/g(X_i)$ , specifically,

$$w(X_i) = \frac{f(X_i)/g(X_i)}{\sum_{k=1}^m f(X_k)/g(X_k)}.$$

- If  $f = Cq$  for some unknown  $C > 0$  but known  $q$ , the same  $w(X_i)$  can be computed with  $f(X_i)$  being replaced by  $q(X_i)$ .

## SIR

Fix  $m$  and  $n$ .

- ① Sample  $X_1, \dots, X_m$  iid from  $g$ .
  - ② Calculate  $w(X_1), \dots, w(X_m)$ .
  - ③ Draw  $n$  draws with replacement from  $X_1, \dots, X_m$  with probabilities  $w(X_1), \dots, w(X_m)$ .
- The returned values  $Y_1, \dots, Y_n$  consist a random sample *approximately* following the density  $f(x)$ .
  - In order for  $Y_1, \dots, Y_n$  to behave like an iid sample from  $f(x)$ , there should be  $m \gg 1$  and  $n/m \rightarrow 0$  as  $m \rightarrow \infty$ .

Let  $h(x)$  be a function proportional to  $f(x)/g(x)$ . The above SIR is a procedure to sample from

$$f(x) \propto h(x)g(x).$$

It can be extended to

$$f(x) \propto \sum_{k=1}^K h_k(x)g_k(x)$$

where  $g_k(x)$  are probability densities and  $h_k(x) \geq 0$  positive functions. However, here the complete form of  $g_k(x)$  has to be known. The procedure can be modified as follows.

## SIR

Fix  $m$  and  $n$ .

- 1 For  $k = 1, \dots, K$ , sample  $X_{k1}, \dots, X_{km}$  iid from  $g_k$ .
- 2 Calculate the weights

$$w_{kj} \propto h_k(X_{kj})$$

for  $k = 1, \dots, K$ ,  $j = 1, \dots, m$ , such that their total is 1.

- 3 Make  $n$  draws with replacement from  $\{X_{kj}\}$  with probabilities  $w_{kj}$ .



To justify SIR, we show that when  $m \gg 1$ ,  $Y_1$  (or  $Y_2$ ,  $Y_3$ , etc) has a distribution approximately equal to  $f$ . Let

$$C = \sum_{k=1}^K \int h_k(x) g_k(x) \mathrm{d}x.$$

Then

$$f(x) = C^{-1} \sum_{k=1}^K h_k(x) g_k(x).$$

We need to show that, if  $A$  is an arbitrary region, then, as  $m \rightarrow \infty$ ,

$$\mathbb{P}\{Y_1 \in A\} \rightarrow \int_A f(x) \mathrm{d}x.$$

Denote by  $\mathbf{X}$  the set of  $X_{kj}$ ,  $k = 1, \dots, K$ ,  $j = 1, \dots, m$ . Given  $\mathbf{X}$ , in the resampling step, the probability that  $X_{kj}$  is drawn to be  $Y_1$  is

$$w_{kj} = \frac{h_k(X_{kj})}{\sum_{k=1}^K \sum_{j=1}^m h_k(X_{kj})}.$$

Then

$$\begin{aligned} \mathbb{P}\{Y_1 \in A \mid \mathbf{X}\} &= \sum_{k=1}^K \sum_{j=1}^m \mathbb{P}\{X_{kj} \text{ is drawn to be } Y_1\} \mathbf{1}\{X_{kj} \in A\} \\ &= \frac{\frac{1}{m} \sum_{k=1}^K \sum_{j=1}^m \mathbf{1}\{X_{kj} \in A\} h_k(X_{kj})}{\frac{1}{m} \sum_{k=1}^K \sum_{j=1}^m h_k(X_{kj})}. \end{aligned}$$

By the Strong Law of Large Numbers, with probability 1, as  $m \rightarrow \infty$ , for each  $k$

$$\frac{1}{m} \sum_{j=1}^m \mathbf{1}\{X_{kj} \in A\} h_k(X_{kj}) \rightarrow \int \mathbf{1}\{x \in A\} h_k(x) g_k(x) dx$$

and so

$$\begin{aligned} & \frac{1}{m} \sum_{k=1}^K \sum_{j=1}^m \mathbf{1}\{X_{kj} \in A\} h_k(X_{kj}) \\ & \rightarrow \int \mathbf{1}\{x \in A\} \sum_{k=1}^K h_k(x) g_k(x) dx = C \int \mathbf{1}\{x \in A\} f(x) dx \end{aligned}$$

In particular, let  $A$  be the entire domain of possible values of  $X_{jk}$ . Then

$$\frac{1}{m} \sum_{k=1}^K \sum_{j=1}^m h_k(X_{kj}) \rightarrow C.$$

So with probability one, as  $m \rightarrow \infty$ ,

$$\mathbb{P}\{Y_1 \in A \mid \mathbf{X}\} \rightarrow \int_A f(x) dx.$$

Since we are only interested in the distribution of  $Y_1$ , we need to remove the dependence on  $\mathbf{X}$ . This requires finding the limit of  $\mathbb{P}\{Y_1 \in A\}$  as  $m \rightarrow \infty$ . Note

$$\mathbb{P}\{Y_1 \in A\} = \mathbb{E}[\mathbb{P}\{Y_1 \in A \mid X_1, \dots, X_m\}].$$

By the Dominated Convergence Theorem,

$$\lim_{m \rightarrow \infty} \mathbb{P}\{Y_1 \in A\} = \mathbb{E}\left[\lim_{m \rightarrow \infty} \mathbb{P}\{Y_1 \in A \mid \mathbf{X}\}\right] = \mathbb{E}\left[\int_A f(x) dx\right] = \int_A f(x) dx$$

# Sampling using Markov chains

The sampling using Markov chains is a very important method.

- Highly adaptive to various difficult sampling problems, especially high-dimensional problems.
- Good performance in most applications despite being approximate.
- Ease of its implementation.

The method is often used to estimate  $\mathbb{E}[h(X)]$  for a function  $h(X)$  of random variable  $X$ . In this context, the method is known as Markov chain Monte Carlo (MCMC). Here we will use this term, but only consider the sampling aspect.

Let  $\mathcal{X}$  be a discrete or Euclidean space. Recall that a stochastic process  $\{X_t, t = 0, 1, 2, \dots\}$  taking values  $\mathcal{X}$  is called a Markov chain if

$$\mathbb{P}\{X_t \in A \mid X_{t-1}, \dots, X_1, X_0\} = \mathbb{P}\{X_t \in A \mid X_{t-1}\} \text{ a.s.}$$

for all  $t > 0$  and (measurable)  $A \subset \mathcal{X}$ . The marginal distribution of  $X_0$  is called the initial distribution of the Markov chain  $(X_t)$ . If every  $X_t$  has a conditional density given  $X_s$ ,  $0 \leq s < t$ , then the condition can be written as

$$p(x_t \mid x_{t-1}, \dots, x_1, x_0) = p(x_t \mid x_{t-1}).$$

The chain is homogeneous if for any  $a$  and  $b$ ,

$$p(x_t = b \mid x_{t-1} = a)$$

is the same for all  $t > 0$ . Denote by  $q(b \mid a)$  to be the conditional density and call  $q(\cdot \mid \cdot)$  the transition kernel of the Markov chain.

A homogeneous Markov chain  $(X_t, t = 0, 1, 2, \dots)$  in general is not stationary, because the density  $p_t$  of  $X_t$  is given by

$$p_t(x) = \int q(x | y) p_{t-1}(y) \, dy$$

and hence is not necessarily the same as  $p_{t-1}$ . However, if  $X_0$  follows a distribution  $\pi$  that satisfies

$$\pi(x) = \int q(x | y) \pi(y) \, dy,$$

then  $p_t \equiv \pi$  for all  $t \geq 0$ . This then implies the stationarity of  $(X_t)$ . Thus  $\pi$  is called a *stationary* distribution of  $(X_t)$ .

Importantly, under certain conditions, no matter the initial distribution  $p_0$ , as  $t \rightarrow \infty$ ,  $p_t \rightarrow \pi$ . The existence and uniqueness of stationary distribution, and the convergence of  $p_t$  are true in a large number of cases, in particular, when  $\mathcal{X}$  is finite and  $p(y | x) > 0$  for all  $x, y \in \mathcal{X}$ .

# Idea of MCMC

- Construct a Markov chain  $(X_t)$  such that it has the target distribution  $f$  as the unique stationary distribution.
- If this can be done, then one has a good reason to expect that  $X_t$  is a random sample from a distribution very close to  $f$ , as long as  $n$  is large enough.
- Useful when it is relatively easy to specify the transition kernel of a Markov chain, however, it is much harder to determine the corresponding stationary distribution(s).

Basis question: how to choose  $q(\cdot | \cdot)$  for  $(X_t)$ , such that  $f$  is a stationary distribution of  $(X_t)$ ?

- If this question is answered, then one can proceed to find conditions to make sure that  $(X_t)$  has a unique stationary distribution and  $p_t$  converges to it.



Suppose  $f$  is known up to a multiplicative factor, i.e

$$f(x) = h(x)/C,$$

where  $C > 0$  is a constant that is possibly intractable. If  $q(\cdot | \cdot)$  satisfies the following *detailed balance* condition,

$$h(x)q(y | x) = h(y)q(x | y) \quad \text{for all } x \neq y \in \mathcal{X},$$

then  $f$  is a stationary distribution of  $(X_t)$ .

**Proof.**

From the condition,  $f(x)q(y | x) = f(y)q(x | y)$  for all  $x$  and  $y \in \mathcal{X}$ . Then

$$\int f(x)q(y | x) dx = \int f(y)q(x | y) dx = f(y) \int q(x | y) dx = f(y)$$

where the last equality is due to the fact that given  $y$ ,  $q(x | y)$  is a probability density. □

## Metropolis-Hastings (MH) algorithm

Detailed balance is the guiding principle for many MCMC algorithms. However, it does not say how to construct a kernel. One answer is provided by the well-known MH algorithm. It uses a *proposal distribution*  $k(y|x)$  to implicitly construct a transition kernel satisfying the detailed balance condition. The construction has some similarity to rejection sampling.

### MH Algorithm

Set an arbitrary  $x_0$  with  $h(x_0) > 0$ . Then, for  $t \geq 0$ , given  $x_t$ , sample  $x_{t+1}$  as follows.

- 1 Draw  $x^* \sim k(\cdot | x_t)$  and  $U \sim \text{Unif}(0, 1)$ .
- 2 Compute the *MH ratio*  $R(x_t, x^*)$ , where

$$R(x, y) = \frac{h(y)k(x|y)}{h(x)k(y|x)}$$

- 3 If  $U \leq \min\{R(x_t, x^*), 1\}$ , then set  $x_{t+1} = x^*$ , else set  $x_{t+1} = x_t$ .

We shall show that the MH algorithm satisfies the detailed balance condition. But first, note that  $R(x_t, x^*)$  is always well-defined, because  $x^*$  can be sampled only if  $f(x_t) > 0$  and  $k(x^* | x) > 0$ .

Since each  $x_{t+1}$  is sampled solely based on  $x_t$ , then  $(x_t)$  is a Markov chain. Let  $q(y | x)$  be its transition kernel. Observe that for  $y \neq x$ ,

$$\begin{aligned} h(x)q(y | x) > 0 &\iff h(x) > 0 \text{ and } q(y | x) > 0 \\ &\iff h(x) > 0, k(y | x) > 0, h(y) > 0, k(x | y) > 0 \\ &\iff h(y) > 0 \text{ and } q(x | y) > 0 \\ &\iff h(y)q(y | x) > 0. \end{aligned}$$

Therefore, the detailed balance condition is satisfied if  $h(x)q(y | x) = 0$ .

Next, if  $h(x)q(y|x) > 0$ , then  $h(y)q(x|y) > 0$ . In this case,

$$\begin{aligned} q(y|x) &= k(y|x) \min\{R(x,y), 1\} \\ &= k(y|x) \min\left\{\frac{h(y)k(x|y)}{h(x)k(y|x)}, 1\right\} \end{aligned}$$

Then

$$\begin{aligned} h(x)q(y|x) &= h(x)k(y|x) \min\left\{\frac{h(y)k(x|y)}{h(x)k(y|x)}, 1\right\} \\ &= \min\{h(y)k(x|y), h(x)k(y|x)\}. \end{aligned}$$

Likewise, since  $h(y)q(x|y) > 0$ , it is also equal to the right hand side of the equation. Therefore, the detailed balance is also satisfied for  $x \neq y$  with  $h(x)k(y|x) > 0$ .

## Bayesian inference

MCMC methods like the MH algorithm are particularly popular tools for Bayesian inference. Let  $f(x|\theta)$  be a parametric family of probability densities and  $p(\theta)$  be a prior distribution of parameter  $\theta$ . Given data  $x$ , the posterior distribution of  $\theta$  is

$$p(\theta|x) \propto h(\theta) := f(x|\theta)p(\theta).$$

A simple strategy is to use the prior  $p(\theta)$  as a proposal distribution; i.e., for any  $\theta_t$ ,  $k(\theta^*|\theta_t) = p(\theta^*)$ . Then the MH ratio is simply the likelihood ratio

$$R(\theta_t, \theta^*) = \frac{h(\theta^*)p(\theta_t)}{h(\theta_t)p(\theta^*)} = \frac{f(x|\theta^*)}{f(x|\theta_t)}.$$

and the M-H algorithm becomes

- Draw  $\theta^* \sim p$  and  $U \sim \text{Unif}(0, 1)$ . If  $f(x|\theta^*) \geq Uf(x|\theta_t)$ , then set  $x_{t+1} = x^*$ ; else set  $x_{t+1} = x_t$ .

# Gibbs sampling

Gibbs sampling is specifically adapted for multidimensional target distributions. It works by sequentially sampling from the *conditional* distributions of parts of the random vector, which are often available in closed form.

Let  $\mathbf{X}$  be a random vector. Suppose we can partition it into random subvectors  $\mathbf{X}_1, \dots, \mathbf{X}_p$ , such that, for each  $i = 1, \dots, p$ , the conditional distribution of  $\mathbf{X}_i$  given all the other  $\mathbf{X}_j$ 's is easy to sample. With a little abuse of notation, denote by

$$f(\mathbf{x}_i \mid \mathbf{x}_1, \dots, \mathbf{x}_{i-1}, \mathbf{x}_{i+1}, \dots, \mathbf{x}_p)$$

the conditional density of  $\mathbf{X}_i$  given  $\mathbf{X}_j = \mathbf{x}_j$  for  $j \neq i$ .

## Gibbs sampler

Set  $t = 0$  and  $\mathbf{x}^{(0)} = (x_1^{(0)}, \dots, x_p^{(0)})$ . Then, for  $t \geq 0$ , given  $\mathbf{x}^{(t)}$ , sample  $\mathbf{x}^{(t+1)}$  as follows.

- 1 Sample  $i$  uniformly from  $\{1, \dots, p\}$ .
- 2 Draw  $x_i^{(t+1)}$  from  $f(x_i | x_1^{(t)}, \dots, x_{i-1}^{(t)}, x_{i+1}^{(t)}, \dots, x_p^{(t)})$
- 3 Set  $x_j^{(t+1)} = x_j^{(t)}$  for all  $j \neq i$ .

Let  $q(\mathbf{y} | \mathbf{x})$  be the transition kernel of the Markov chain in the Gibbs sampler. Note that for  $\mathbf{x} \neq \mathbf{y}$ ,  $q(\mathbf{y} | \mathbf{x}) > 0$  only when there is exactly one  $i$  such that  $\mathbf{x}_i \neq \mathbf{y}_i$  and for all  $j \neq i$ ,  $\mathbf{x}_j = \mathbf{y}_j$ . In this case, writing  $\mathbf{x}_{-i} = (\mathbf{x}_1, \dots, \mathbf{x}_{i-1}, \mathbf{x}_{i+1}, \dots, \mathbf{x}_p)$ ,

$$\begin{aligned} f(\mathbf{x})q(\mathbf{y} | \mathbf{x}) &= \frac{1}{n}f(\mathbf{x}_1, \dots, \mathbf{x}_p)f(\mathbf{y}_i | \mathbf{x}_{-i}) \\ &= \frac{1}{n}f(\mathbf{x}_i | \mathbf{x}_{-i})f(\mathbf{x}_{-i})f(\mathbf{y}_i | \mathbf{x}_{-i}) \end{aligned}$$

where the factor  $1/n$  is the probability that  $i$  is sampled. Because,  $\mathbf{y}_{-i} = \mathbf{x}_{-i}$ ,

$$\begin{aligned} f(\mathbf{x})q(\mathbf{y} | \mathbf{x}) &= \frac{1}{n}f(\mathbf{x}_i | \mathbf{y}_{-i})f(\mathbf{y}_{-i})f(\mathbf{y}_i | \mathbf{y}_{-i}) \\ &= \frac{1}{n}f(\mathbf{x}_i | \mathbf{y}_{-i})f(\mathbf{y}) \\ &= f(\mathbf{y})q(\mathbf{x} | \mathbf{y}). \end{aligned}$$

Therefore, the detailed balance condition is satisfied.



Usually, instead of randomly sampling the index  $i$ , the following variant is used. Note that for  $p > 1$ , the variant does not satisfy the detailed balance condition. However, under certain conditions, the Markov chain  $\mathbf{x}^{(t)}$  generated by the sampler still has  $f$  as the stationary distribution.

### Gibbs sampler with cycles

Set  $t = 0$  and  $\mathbf{x}^{(0)} = (x_1^{(0)}, \dots, x_p^{(0)})$ . Then, for  $t \geq 0$ , given  $\mathbf{x}^{(t)}$ , sample  $\mathbf{x}^{(t+1)}$  as follows.

Draw  $x_1^{(t+1)}$  from  $f(x_1 | x_2^{(t)}, x_3^{(t)}, x_4^{(t)}, \dots, x_p^{(t)})$ .

Draw  $x_2^{(t+1)}$  from  $f(x_2 | x_1^{(t+1)}, x_3^{(t)}, x_4^{(t)}, \dots, x_p^{(t)})$ .

Draw  $x_3^{(t+1)}$  from  $f(x_3 | x_1^{(t+1)}, x_2^{(t+1)}, x_4^{(t)}, \dots, x_p^{(t)})$ .

.....

Draw  $x_{p-1}^{(t+1)}$  from  $f(x_{p-1} | x_1^{(t+1)}, x_2^{(t+1)}, \dots, x_{p-2}^{(t+1)}, x_p^{(t)})$ .

Draw  $x_p^{(t+1)}$  from  $f(x_p | x_1^{(t+1)}, x_2^{(t+1)}, \dots, x_{p-1}^{(t+1)})$ .

# Simulating Brownian motions

A *standard* one-dimensional (1D) Brownian motion (BM) on  $[0, \infty)$  is a stochastic process  $\{W(t) : t \geq 0\}$ , such that

- 1  $W(0) = 0$ ;
- 2  $W(t)$  is a continuous function with probability one;
- 3 for any  $0 \leq t_0 < t_1 < \cdots < t_k < \cdots$ , the increments

$$W(t_i) - W(t_{i-1}), \quad i = 1, 2, \dots$$

are independent of each other; and

- 4 for any  $0 \leq s < t$ ,  $W(t) - W(s) \sim N(0, t - s)$ .

For constants  $\mu, \sigma^2 > 0$ , if  $\sigma^{-1}[X(t) - \mu t]$  is a standard BM, then  $X$  is called a BM with drift  $\mu$  and diffusion coefficient  $\sigma^2$ , denoted  $X \sim \text{BM}(\mu, \sigma^2)$ . More generally, a BM with deterministic but time-varying drift  $\mu(t)$  and diffusion coefficient  $\sigma^2(t)$  satisfies the SDE

$$dX(t) = \mu(t) dt + \sigma(t) dW(t),$$

i.e.

$$X(t) = X(0) + \int_0^t \mu(s) ds + \int_0^t \sigma(s) dW(s).$$

It has a continuous sample path almost surely and independent increments. For  $0 \leq s < t < \infty$ ,  $X(t) - X(s)$  is normally distributed with

$$\mathbb{E}[X(t) - X(s)] = \int_s^t \mu(x) dx,$$

$$\text{Var}[X(t) - X(s)] = \int_s^t \sigma^2(x) dx.$$

## Random walk construction

Suppose the value of  $X(0)$  is either known and fixed or can be sampled.

Given  $0 = t_0 < t_1 < t_2 < \dots < t_n$ ,

$$D_i = X(t_{i+1}) - X(t_i), \quad i = 0, \dots, n-1,$$

are independent and can be written as

$$D_i = \int_{t_i}^{t_{i+1}} \mu(x) dx + \sqrt{\int_{t_i}^{t_{i+1}} \sigma^2(x) dx} Z_{i+1},$$

with  $Z_1, \dots, Z_n$  iid  $\sim N(0, 1)$ . This yields the following algorithm to sample  $X(t_1), \dots, X(t_n)$ .

- ① Set or sample  $X(0)$
- ② Sample  $Z_1, \dots, Z_n$  iid  $\sim N(0, 1)$
- ③ For  $i = 0, \dots, n-1$ , compute  $D_i$  and set  $X(t_{i+1}) = X(t_i) + D_i$ .
- ④ Return  $(X(0), X(t_1), \dots, X(t_n))$

# Brownian bridge construction

For BM  $W(t)$  with drift  $\mu(t)$  and diffusion coefficient  $\sigma^2(t)$ , the random walk construction generates  $W(t_i)$  from left to right. Alternatively, one may first generate the final value  $W(t_n)$ , then sample  $W(t_{\lfloor n/2 \rfloor})$  conditional on  $W(t_n)$  and progressively filling in intermediate values.

- 1 Useful in implementing variance reduction techniques and low-discrepancy methods;
- 2 Allows the “time resolution” of the sample to be increased based on already sampled  $W(t_i)$ ;
- 3 The key is the conditioning formula for jointly normal random vectors.

Recall the following conditioning formula.

Let  $\mathbf{X}_1 \in \mathbb{R}^m$ ,  $\mathbf{X}_2 \in \mathbb{R}^n$  be jointly normal, such that

$$\mathbb{E}\mathbf{X}_i = \boldsymbol{\mu}_i, \quad \text{Cov}(\mathbf{X}_i, \mathbf{X}_j) = \boldsymbol{\Sigma}_{ij}.$$

If  $\boldsymbol{\Sigma}_{22}$  is of full rank, then, conditional on  $\mathbf{X}_2 = \mathbf{x}$ ,  $\mathbf{X}_1 \sim N(\boldsymbol{\mu}_c, \boldsymbol{\Sigma}_c)$ , where

$$\boldsymbol{\mu}_c = \boldsymbol{\mu}_1 + \boldsymbol{\Sigma}_{12}\boldsymbol{\Sigma}_{22}^{-1}(\mathbf{x} - \boldsymbol{\mu}_2), \quad \boldsymbol{\Sigma}_c = \boldsymbol{\Sigma}_{11} - \boldsymbol{\Sigma}_{12}\boldsymbol{\Sigma}_{22}^{-1}\boldsymbol{\Sigma}_{21}.$$

Note that  $\boldsymbol{\Sigma}_c$  is independent of  $\mathbf{x}$ .

We shall focus the Brownian bridge construction for  $\text{BM}(0, 1)$ . The more general case follows the same idea but a bit more complex.

To start with, given  $0 \leq a < t < b$ , consider how to sample  $W(t)$  given  $W(a) = x$  and  $W(b) = y$ :

$$\begin{array}{ccc} a & t & b \\ W(a) = x & W(t) = ? & W(b) = y. \end{array}$$

Since  $\text{Cov}(W(t_1), W(t_2)) = t_1$  for any  $0 \leq t_1 < t_2$ ,

$$\begin{pmatrix} W(t) \\ W(a) \\ W(b) \end{pmatrix} \sim N \left( 0, \begin{pmatrix} t & a & t \\ a & a & a \\ t & a & b \end{pmatrix} \right)$$

Take  $\mathbf{X}_1 = W(t)$ ,  $\mathbf{X}_2 = (W(a), W(b))'$ . Then  $\mu_1 = 0$ ,  $\mu_2 = \mathbf{0}$ , and

$$\Sigma_{11} = t, \quad \Sigma_{22} = \begin{pmatrix} a & a \\ a & b \end{pmatrix}, \quad \Sigma_{12} = (a, t).$$

Then, conditioning on  $W(a) = x$ ,  $W(b) = y$ ,  $W(t) \sim N(\mu, \sigma^2)$ , with

$$\begin{aligned}\mu &= \mu_1 + \Sigma_{12}\Sigma_{22}^{-1}(\mathbf{x} - \mu_2) = (a, t) \begin{pmatrix} a & a \\ a & b \end{pmatrix}^{-1} \begin{pmatrix} x \\ y \end{pmatrix} \\ &= \frac{(b-t)x + (t-a)y}{b-a},\end{aligned}$$

and

$$\sigma^2 = \Sigma_{11} - \Sigma_{12}\Sigma_{22}^{-1}\Sigma_{21} = t - (a, t) \begin{pmatrix} a & a \\ a & b \end{pmatrix}^{-1} \begin{pmatrix} a \\ t \end{pmatrix} = \frac{(t-a)(b-t)}{b-a}.$$

Note the conditional mean of  $W(t)$  is the linear interpolation of  $x$  and  $y$ .



More generally, suppose we have  $W(s_i)$ , for  $0 = s_0 < s_1 < \dots < s_n$  and wish to insert more time points. There are two important facts that allow this to be done. First, if  $s_i < t < s_{i+1}$ , then

$$\begin{aligned} &\text{Conditional distribution of } W(t) \text{ given } \{W(s_j), j = 1, \dots, n\} \\ &= \text{Conditional distribution of } W(t) \text{ given } \{W(s_i), W(s_{i+1})\}. \end{aligned}$$

Second, for each interval  $(s_{i-1}, s_i)$ , if we choose exactly one  $t_i$ , then  $W(t_1), \dots, W(t_n)$  are *conditionally independent*:

$$\begin{aligned} &\text{given } W(s_j) = x_j, j = 1, \dots, n \\ &(W(t_1), \dots, W(t_n)) \sim (\mu_1 + \sigma_1 Z_1, \dots, \mu_n + \sigma_n Z_n) \end{aligned}$$

where  $Z_1, \dots, Z_n$  are iid  $\sim N(0, 1)$ , and

$$\mu_i = \frac{(s_i - t_i)x_{i-1} + (t_i - s_{i-1})x_i}{s_i - s_{i-1}}, \quad \sigma_i^2 = \frac{(t_i - s_{i-1})(s_i - t_i)}{s_i - s_{i-1}}$$

For general

$$W(t) \sim \text{BM}(\mu(t), \sigma^2(t))$$

one has to make several changes. First, define

$$u(t) = \int_0^t \mu(x) \, dx, \quad d^2(t) = \int_0^t \sigma^2(x) \, dx$$

Then for  $a < t < b$ , use

$$\begin{aligned} \boldsymbol{\mu}_1 &= u(t), \quad \boldsymbol{\mu}_2 = \begin{pmatrix} u(a) \\ u(b) \end{pmatrix}, \\ \boldsymbol{\Sigma}_{11} &= d^2(t), \quad \boldsymbol{\Sigma}_{22} = \begin{pmatrix} d^2(a) & d^2(a) \\ d^2(a) & d^2(b) \end{pmatrix} \\ \boldsymbol{\Sigma}_{12} &= (d^2(a), d^2(t)). \end{aligned}$$

in the conditioning formula.

## Geometric BM: one dimension, fixed drift

If  $\ln S(t)$  is a BM, then  $S(t)$  is called a *geometric BM*. It is often specified by an SDE. For example, if  $S$  satisfies the SDE

$$\frac{dS(t)}{S(t)} = \mu dt + \sigma dW(t),$$

then we denote  $S \sim \text{GBM}(\mu, \sigma^2)$ . To solve the SDE, we may guess

$$S(t) = S(0)e^{at+bW(t)} = f(t, W(t)),$$

where  $f(t, x) = S(0)e^{at+bx}$ . Then by Itô's formula,

$$\begin{aligned} dS(t) &= f'_t dt + f'_x dW(t) + \frac{1}{2} f''_{xx} (dW(t))^2|_{x=W(t)} \\ &= aS(t) dt + bS(t) dW(t) + \frac{b^2}{2} S(t) dt. \end{aligned}$$

Comparing the SDEs,  $b = \sigma$ ,  $a = \mu - \sigma^2/2$ , and so

$$S(t) = S(0)e^{(\mu - \sigma^2/2)t + \sigma W(t)}.$$

## Gaussian short rate models

The instantaneous continuously compounded short rate at time  $t$ , denoted by  $r(t)$ , is important. First, it determines the value  $\beta(t)$  of a risk-free asset at time  $t$  by

$$\beta(t) = \beta(0) \exp \left( \int_0^t r(u) \, du \right).$$

Second, letting

$B(t, T)$  = time- $t$  price of a bond paying 1 at  $T$ ,

where  $T > 0$  is fixed and  $t \in [0, T]$ , from the formula of  $\beta(t)$ , it is seen

$$B(0, T) = \mathbb{E} \left[ \exp \left( - \int_0^T r(u) \, du \right) \right], \quad B(t, T) = \mathbb{E} \left[ \exp \left( - \int_t^T r(u) \, du \right) \right]$$

and the instantaneous continuously compounded forward rate fixed at  $t$  for  $T > t$  is

$$f(t, T) = - \frac{\partial}{\partial T} \ln B(t, T).$$

It is thus important to model  $r(t)$ . Two simple but important models are as follows.

*Vasicek model* is

$$dr(t) = \alpha(b(t) - r(t))dt + \sigma dW(t)$$

- ①  $\alpha$  and  $\sigma$ : positive deterministic constants
- ②  $b(t) > 0$ : deterministic, specified according to bond prices

Under Vasicek model  $r(t)$  is pulled toward  $b(t)$  at speed  $\alpha > 0$ . Long-run interest rate can be modeled by setting  $b(t) \equiv b$ . The result is the Ornstein-Uhlenbeck process

*Continuous-time Ho-Lee model* is

$$dr(t) = g(t)dt + \sigma dW(t).$$

In practice,  $g(t)$  may be specified according to observed bond prices.

Both models are special cases to the *Gaussian Markov process*

$$dr(t) = [g(t) + h(t)r(t)] dt + \sigma(t) dW(t)$$

where  $g$ ,  $h$  and  $\sigma$  are deterministic functions of  $t$ . Let

$$H(t) = \int_0^t h(s) ds.$$

By Itô's formula

$$d(e^{-H(t)}r(t)) = e^{-H(t)}g(t) dt + e^{-H(t)}\sigma(t) dW(t),$$

and so the general solution is

$$r(t) = e^{H(t)}r(0) + \int_0^t e^{H(t)-H(s)}g(s) ds + \int_0^t e^{H(t)-H(s)}\sigma(s) dW(s)$$

In particular, Vasicek model has solution

$$r(t) = e^{-\alpha t} r(0) + \alpha \int_0^t e^{-\alpha(t-s)} b(s) \, ds + \sigma \int_0^t e^{-\alpha(t-s)} \, dW(s)$$

and Ho-Lee model has solution

$$r(t) = r(0) + \int_0^t g(s) \, ds + \sigma W(t).$$

To simulate the models, first consider the general process. Its simulation is based on the following Markov property: conditional on  $r(\tau)$ ,  $\tau \in [0, u]$ , for  $t > u$ ,

$$r(t) \sim N \left( e^{H(t)-H(u)} r(u) + \mu(u, t), \sigma_r^2(u, t) \right),$$

with

$$\mu(u, t) = \int_u^t e^{H(t)-H(s)} g(s) \, ds, \quad \sigma_r^2(u, t) = \int_u^t e^{2[H(t)-H(s)]} \sigma^2(s) \, ds.$$

This can be seen by using

$$e^{-H(t)} r(t) = e^{-H(u)} r(u) + \int_u^t e^{-H(s)} g(s) \, ds + \int_u^t e^{-H(s)} \sigma(s) \, dW(s).$$



This result implies that to sample  $r(t)$  and  $0 = t_0 < t_1 < t_2 < \dots$ , one can do the following.

- 1 Sample  $Z_1, Z_2, \dots$  iid  $\sim N(0, 1)$ .
- 2 For  $i = 0, 1, \dots$ , set

$$r(t_{i+1}) = e^{H(t_{i+1}) - H(t_i)} r(t_i) + \mu(t_i, t_{i+1}) + \sigma_r(t_i, t_{i+1}) Z_{i+1}$$

The above procedure can be applied to Vasicek model and Ho-Lee model easily. Just note that in Vasicek model,  $H(t) = -\alpha t$  and

$$\mu(u, t) = \alpha \int_u^t e^{-\alpha(t-s)} b(s) \, ds, \quad \sigma_r^2(u, t) = \frac{\sigma^2}{2\alpha} (1 - e^{-2\alpha(t-u)}).$$

while in Ho-Lee model,  $H(t) = 0$  and

$$\mu(u, t) = \int_u^t g(s) \, ds, \quad \sigma_r^2(u, t) = (t - u) \sigma^2.$$

We need to know  $b(s)$  and  $g(s)$  in order to simulate the models. The question here is whether they can be set according to observations. In particular, since the forward rates  $f(0, t)$  are observed, it would be nice if  $b$  and  $g$  can be set by it.

It turns out this indeed is the case. To get formulas for  $b$  and  $g$  in terms of  $f(0, t)$ , for the general case,  $B(0, T) = \mathbb{E}[e^{-R(T)}]$ , with  $R(T) = \int_0^T r(u) du \sim \text{Normal}$ . Then

$$B(0, T) = \exp \left( -\mathbb{E}(R(T)) + \frac{1}{2} \text{Var}(R(T)) \right).$$

However,

$$B(0, T) = \exp \left( - \int_0^T f(0, t) dt \right).$$

Therefore

$$\mathbb{E}(R(T)) - \frac{1}{2} \text{Var}(R(T)) = \int_0^T f(0, t) dt.$$

This relation can be used to set  $b(s)$  or  $g(s)$  based on observed  $f(0, t)$ .

For Ho-Lee model, this is particularly nice. One get

$$g(t) = \frac{\partial f(0, t)}{\partial t} + \sigma^2 t.$$

Then, in the simulation, one can use

$$\begin{aligned} r(t_{i+1}) &= r(t_i) + \int_{t_i}^{t_{i+1}} g(s) \, ds + \sigma \sqrt{t_{i+1} - t_i} Z_{i+1} \\ &= r(t_i) + [f(0, t_{i+1}) - f(0, t_i)] + \frac{1}{2} \sigma^2 [t_{i+1}^2 - t_i^2] + \sigma \sqrt{t_{i+1} - t_i} Z_i. \end{aligned}$$

In the iteration, only  $f(0, t_i)$  are needed (and can be observed), and there is no need to differentiate  $f(0, t)$ .

For Vasicek model,

$$r(t_{i+1}) = e^{-\alpha(t_{i+1}-t_i)}r(t_i) + \mu(t_i, t_{i+1}) + \sigma\sqrt{\frac{1 - e^{-2\alpha(t_{i+1}-t_i)}}{2\alpha}}Z_{i+1}.$$

It turns out

$$\mu(u, t) = [f(0, t) + I(0, t)] - e^{-\alpha(t-u)}[f(0, u) + I(0, u)].$$

where

$$I(u, t) = \frac{\sigma^2}{2\alpha^2}(1 - e^{-\alpha(t-u)})^2$$

is computable. The above iteration thus is determined by  $f(0, t)$ , which is observable.

## Square-root diffusions

The process satisfies the SDE

$$dr(t) = \alpha(b - r(t))dt + \sigma\sqrt{r(t)}dW(t).$$

where  $W \sim \text{BM}(0, 1)$ , and  $\alpha, b > 0$ . It was proposed by Cox, Ingersoll and Ross (CIR) as a model of the short rate. In the model,  $r(t)$  is pulled towards  $b$  at rate  $\alpha$ , and as  $r(t) \downarrow 0$ ,  $\sigma\sqrt{r(t)} \rightarrow 0$ , results in  $r(t)$  always being positive.

The model also appears in volatility models. For example

$$\begin{aligned}\frac{dS(t)}{S(t)} &= \mu dt + \sqrt{V(t)}dW_1(t) \\ dV(t) &= \alpha(b - V(t))dt + \sigma\sqrt{V(t)}dW_2(t),\end{aligned}$$

with  $W_1, W_2 \text{ iid } \sim \text{BM}(0, 1)$ .

Euler discretization can be used to simulate the model *approximately*

$$r(t_{i+1}) \approx r(t_i) + \alpha(b - r(t_i))(t_{i+1} - t_i) + \sigma\sqrt{r(t_i)^+(t_{i+1} - t_i)} Z_{i+1}$$

with  $Z_1, Z_2, \dots$  iid  $\sim N(0, 1)$ .

For *exact* simulation, note that  $r(t)$  is Markov. Conditional on  $r(t)$ ,

$$\frac{r(t + \Delta t)}{C(\Delta t)} \sim \chi_d^2 \left( \frac{e^{-\alpha\Delta t} r(t)}{C(\Delta t)} \right), \quad \Delta t > 0,$$

where

$$C(\Delta t) = \frac{\sigma^2(1 - e^{-\alpha\Delta t})}{4\alpha}, \quad d = \frac{4\alpha b}{\sigma^2},$$

and  $\chi_d^2(\lambda)$  is the non-central Gamma distribution with degree of freedom  $d$  and noncentrality parameter  $\lambda \geq 0$ .

The distribution  $\chi_d^2(\lambda)$  is rather complicated. However, it can be characterized as follows. For  $d < 1$ ,

$$\chi_d^2(\lambda) = \chi_{d+2N}^2, \quad N \sim \text{Poisson}(\lambda/2);$$

while for  $d > 1$ ,

$$\chi_d^2(\lambda) = (Z + \sqrt{\lambda})^2 + \chi_{d-1}^2, \quad Z \sim N(0, 1).$$

Simulation of CIR on  $0 = t_0 < t_1 < \dots < t_n$

Case 1:  $d > 1$

for  $i = 0, \dots, n-1$

set  $c = C(t_{i+1} - t_i)$ ,  $\lambda = e^{-\alpha(t_{i+1} - t_i)} r(t_i) / c$

sample  $Z \sim N(0, 1)$ ,  $X \sim \chi_{d-1}^2$

set  $r(t_{i+1}) = c[(Z + \sqrt{\lambda})^2 + X]$

Case 2:  $d < 1$

for  $i = 0, \dots, n-1$

set  $c = C(t_{i+1} - t_i)$ ,  $\lambda = e^{-\alpha(t_{i+1} - t_i)} r(t_i) / c$

sample  $N \sim \text{Poisson}(\lambda/2)$ ,  $X \sim \chi_{d+2N}^2$

set  $r(t_{i+1}) = cX$



In the above CIR model,  $\alpha$ ,  $b$  and  $\sigma$  constants. In principle, they can be made deterministic but time-varying. In particular,

$$dr(t) = \alpha(b(t) - r(t)) dt + \sigma\sqrt{r(t)} dW(t)$$

is frequently used to make the bond price function

$$t \rightarrow \mathbb{E} \left[ \exp \left( - \int_0^t r(u) du \right) \right]$$

match a set of observed bond prices  $B(0, t)$ . In this case, Euler discretization can be used to sample  $r(t_i)$  approximately,

$$r(t_{i+1}) \approx r(t_i) + \alpha(b(t_i) - r(t_i))(t_{i+1} - t_i) + \sigma\sqrt{r(t_i)}(t_{i+1} - t_i) Z_{i+1}$$

with  $Z_1, Z_2, \dots$  iid  $\sim N(0, 1)$ .

The CIR can be extended in many ways. A simple multifactor extension is to regard  $r(t)$  as the sum of many factors,

$$r(t) = \sum_{i=1}^n X_i(t),$$

where  $X_i$  are independent square-root diffusion processes

$$dX_i(t) = \alpha_i(b_i - X_i(t)) dt + \sigma_i \sqrt{X_i(t)} dW_i(t)$$

with  $W_1(t), \dots, W_n(t)$  iid  $\sim \text{BM}(0, 1)$ .

Finally, we show that the single factor CIR can be regarded as the sum of squares of independent Ornstein–Uhlenbeck processes. Specially, let  $X_1(t), \dots, X_n(t)$  be independent Ornstein–Uhlenbeck processes

$$dX_i(t) = -\frac{\alpha}{2}X_i(t)dt + \frac{\sigma}{2}dW_i(t)$$

where  $\alpha, \sigma$  are constants and  $W_1(t), \dots, W_n(t)$  are iid  $\sim \text{BM}(0, 1)$ , then

$$Y(t) = \sum_{i=1}^n X_i^2(t)$$

is a square-root diffusion

$$dY(t) = \alpha \left( \frac{\sigma^2 n}{4\alpha} - Y(t) \right) dt + \sigma \sqrt{Y(t)} d\tilde{W}(t),$$

where  $\tilde{W}(t) \sim \text{BM}(0, 1)$  is defined via

$$d\tilde{W}(t) = \sum_{i=1}^n \frac{X_i(t)}{\sqrt{Y(t)}} dW_i(t).$$

Because of the above result, for  $n$  integer-valued, one can use the following alternative to simulate  $r(t)$ ,

$$r(t_{i+1}) = \frac{1}{n} \sum_{j=1}^n \left[ \sqrt{\frac{r(t_i)}{e^{\alpha \Delta t_i}}} + \sqrt{n C(\Delta t_i)} Z_{i+1,j} \right]^2$$

where  $Z_{i+1,j}$  are iid  $\sim N(0, 1)$  and  $\Delta t_i = t_{i+1} - t_i$ .

## Constant elasticity of variance (CEV) process

The process is defined via

$$dS(t) = \mu S(t) dt + \sigma S(t)^{\beta/2} dW(t),$$

or

$$\frac{dS(t)}{S(t)} = \mu dt + \sigma S(t)^{(\beta-2)/2} dW(t).$$

The model introduces dependency of the volatility on the asset price. It turns out that  $X(t) = S(t)^{2-\beta}$  is a square-root diffusion:

$$dX(t) = \left[ \frac{\sigma^2}{2}(2-\beta)(1-\beta) + \mu(2-\beta)X(t) \right] dt + \sigma(2-\beta)\sqrt{X(t)} dW(t),$$

so one can simulate  $S(t)$  by first simulating  $X(t)$  using the procedure for square-root diffusions, and then taking  $S(t) = X(t)^{1/(2-\beta)}$ .

# Processes with Jumps

Stock prices often exhibit statistical properties that cannot be captured by diffusion models. One argument is that this may be due to peculiar sudden events affecting individual assets but not the market as a whole. To take into account of such random events, stochastic processes with jumps can be used.

# Merton's jump-diffusion model

Incorporate jumps: jump times and jump sizes.

Let  $0 < \tau_1 < \tau_2 < \dots$  be the random arrival times of events. Suppose that between the events,  $S(t)$  is a diffusion as usual, e.g.

$$\frac{dS(t)}{S(t)} = \mu dt + \sigma dW(t), \quad t \neq \tau_i,$$

on the other hand, the effect of the  $i$ th event on  $S(t)$  is to suddenly change it by  $Y_i$ -fold

$$S(\tau_i) = S(\tau_i-) Y_i$$

where  $S(\tau_i-) = \lim_{t \rightarrow \tau_i-} S(t)$  is the value of  $S(t)$  just before the  $i$ th event.

To get a single SDE accounting for both diffusion and jumps, let

$$N(t) = (\text{Number of events by time } t) = \# \{i : \tau_i \leq t\}.$$

The process  $N(t)$  is identified with the process of arrival times  $\{\tau_i\}$ . Let

$$J(t) = \sum_{i: \tau_i \leq t} (Y_i - 1) = \sum_{i=1}^{N(t)} (Y_i - 1).$$

The desired single SDE is

$$\frac{dS(t)}{S(t-)} = \mu dt + \sigma dW(t) + dJ(t).$$



To see this, note  $J(t)$  is a step function that only has jumps at arrivals of events,

$$dJ(t) = \begin{cases} 0 & t \notin \{\tau_i\}, \\ Y_i - 1 & t = \tau_i. \end{cases}$$

Between jumps,  $S(t)$  is the usual diffusion,  $dS(t) = S(t)[\mu dt + \sigma dW(t)]$ . As  $dJ(t) = 0$  and  $S$  is continuous at  $t$ , one can write

$$dS(t) = S(t-)[\mu dt + \sigma dW(t) + dJ(t)].$$

On the other hand, at time  $t = \tau_i$ ,

$$dS(t) = S(\tau_i) - S(\tau_i-) = S(t-)(Y_i - 1) = S(t-)dJ(t).$$

Now  $\mu dt + \sigma dW(t)$  is negligible comparing to  $dJ(t)$ . So one can write, again,

$$dS(t) = S(t-)[\mu dt + \sigma dW(t) + dJ(t)].$$

Suppose in the jump-diffusion process,

- 1)  $\mu$  and  $\sigma$  are constants,
- 2)  $W(t) \sim \text{BM}(0, 1)$ ,
- 3)  $(\tau_i, Y_i)$  are independent of  $W(t)$ ;  $\tau_i$  and  $Y_i$  can be dependent;
- 4)  $Y_i > 0, i = 1, 2, \dots$

Then

$$S(t) = S(0)e^{(\mu - \sigma^2/2)t + \sigma W(t)} \prod_{i=1}^{N(t)} Y_i.$$

The process can be decomposed into two factors. One is the geometric BM

$$S(0)e^{(\mu - \sigma^2/2)t + \sigma W(t)}$$

The second one is  $\prod_{i=1}^{N(t)} Y_i$ , which is the compounded jump up to time  $t$ . We need to know how to simulate the second one.

# Simulating jump times

The simplest model for jump times is Poisson process. A Poisson process with rate  $\lambda$  is defined by two properties

- 1) for disjoint time intervals, the happenings of events in them are independent of each other; and
- 2) for  $I = [t, t + \Delta]$ ,

$$P\{\text{no event occurs in } I\} = 1 - \lambda\Delta + o(\Delta)$$

$$P\{\text{one event occurs in } I\} = \lambda\Delta$$

$$P\{\text{two or more events occur in } I\} = o(\Delta)$$

There are two ways to simulate a Poisson process with rate  $\lambda$ . The first one utilizes the fact that the “interarrival times”

$$\tau_{i+1} - \tau_i \quad \text{are i.i.d.} \quad \sim \text{Exp}(\lambda)$$

The second one utilizes the following facts

$$N(t) \sim \text{Poisson}(\lambda t), \quad \text{any } t > 0.$$

and given  $N(t) = n$ ,  $\tau_1, \dots, \tau_n$  are the order statistics of  $X_1, \dots, X_n$  iid  $\sim \text{Unif}(0, t)$ , i.e.,  $\tau_1 = X_{(1)}$ , the smallest among  $X_i$ ,  $\tau_2 = X_{(2)}$ , the second smallest, and so on.

To sample jump times in  $[0, T]$ , the methods are implemented as follows.

### Method 1

```
set  $\tau_0 = 0$   
for  $i = 1, 2 \dots$   
  sample  $X \sim \text{Exp}(1)$  and set  $\tau_i = \tau_{i-1} + \lambda X$   
  if  $\tau_i > T$   
    return  $\tau_1, \dots, \tau_{i-1}$ 
```

### Method 2

```
sample  $N \sim \text{Poisson}(\lambda T)$   
sample  $U_1, \dots, U_N \text{ iid } \sim \text{Unif}[0, 1]$   
return  $TU_{(1)}, \dots, TU_{(N)}$ 
```

An inhomogeneous Poisson process with intensity function  $\lambda(t)$  is a process satisfying

- 1 for disjoint time intervals, the happenings of events in them are independent of each other; and
- 2 for  $I = [t, t + \Delta]$ ,

$$P\{\text{no event occurs in } I\} = 1 - \lambda(t)\Delta + o(\Delta)$$

$$P\{\text{one event occurs in } I\} = \lambda(t)\Delta$$

$$P\{\text{two or more events occur in } I\} = o(\Delta)$$

The process has the following important property. Given  $T > 0$ , let

$$Z = \int_0^T \lambda(t) dt.$$

Note that  $f(t) = \lambda(t)/Z$  is a probability density on  $[0, T]$ . Then

- 1  $N(T) \sim \text{Poisson}(Z)$ ; and
- 2 conditioning on  $N(T) = n$ , the arrival times in time interval  $[0, T]$  are the order statistics of  $X_1, \dots, X_n$  i.i.d.  $\sim f(t)$ .

To sample arrival times in  $[0, T]$ , one can then do the following.

- 1 Sample  $N \sim \text{Poisson}(Z)$ .
- 2 Sample  $X_1, \dots, X_N$  iid  $\sim f(t)$  and sort them into  $X_{(1)} \leq \dots \leq X_{(N)}$ .
- 3 Return  $X_{(1)}, \dots, X_{(N)}$ .

If  $\lambda(t)$  is complicated so that either  $Z$  is unknown or  $f(t)$  is difficult to sample from, then one can use the “thinning method”, which is similar to rejection sampling. Let  $\lambda_0(t)$  be a function satisfying  $\lambda(t) \leq \lambda_0(t)$ . Suppose it is easy to sample a Poisson process with intensity function  $\lambda_0(t)$ .

- 1 Sample  $\{\tau_i\}$  from a Poisson process with intensity  $\lambda_0(t)$  in  $[0, T]$ .
- 2 For each  $i$ , retain  $\tau_i$  with probability  $\lambda(\tau_i)/\lambda_0(\tau_i)$ .
- 3 Return the retained  $\tau_i$



# Simulating a jump-diffusion process at fixed time points

Assume

- ①  $\tau_1, \tau_2, \dots$  form a Poisson process of intensity  $\lambda(t)$
- ②  $Y_1, Y_2, \dots$  are iid
- ③  $\{\tau_1, \tau_2, \dots\}$ ,  $\{Y_1, Y_2, \dots\}$  and  $W(t)$  are mutually independent.

Recall the SDE for  $S(t)$  is

$$\frac{dS(t)}{S(t-)} = \mu dt + \sigma dW(t) + dJ(t),$$

with

$$J(t) = \sum_{i=1}^{N(t)} (Y_i - 1).$$

In general, if  $W_i$  are independent random variables that are also independent of a Poisson process  $N(t)$ , then the process

$$M(t) = \sum_{i=1}^{N(t)} W_i$$

is called a compound Poisson process.

Under the assumptions,  $J(t)$  is a compound Poisson process. Let  $X(t) = \ln S(t)$ . Then

$$X(t) = (\mu - \frac{1}{2}\sigma^2)t + \sigma W(t) + \sum_{j=1}^{N(t)} \ln Y_j,$$

where

$$D(t) = \sum_{j=1}^{N(t)} \ln Y_j$$

is another compound Poisson process.

For  $0 < t < u$ ,

$$X(u) = X(t) + (\mu - \tfrac{1}{2}\sigma^2)(u - t) + \sigma[W(u) - W(t)] + \sum_{j=N(t)+1}^{N(u)} \ln Y_j.$$

It is seen that only the total number of jumps between  $t$  and  $u$  matters in the contribution of the compound Poisson process. The actual jump times are not important. Therefore, given time points  $t_1 < t_2 < \dots$ ,  $X(t_i)$  can be simulated as follows,

- ① Sample  $Z \sim N(0, 1)$
- ② Sample  $N \sim \text{Poisson}(\theta_i)$ , where  $\theta_i = \int_{t_i}^{t_{i+1}} \lambda(t) dt$ .
- ③ If  $N = 0$ , set  $M = 0$ ; otherwise, sample  $\ln Y_1, \dots, \ln Y_N$  and set  $M = \sum_{i=1}^N \ln Y_i$ .
- ④ Set

$$X(t_{i+1}) = X(t_i) + (\mu - \tfrac{1}{2}\sigma^2)(t_{i+1} - t_i) + \sigma\sqrt{t_{i+1} - t_i}Z + M.$$

## Simulating at jump times

If one is interested in what may happen at jump times, then it is necessary to simulate the point process  $\tau_1, \tau_2, \dots$ . Note that, if  $\tau_j < \tau_{j+1}$  are adjacent jump times, then

$$X(\tau_{j+1}) = X(\tau_j) + (\mu - \frac{1}{2}\sigma^2)(\tau_{j+1} - \tau_j) + \sigma[W(\tau_{j+1}) - W(\tau_j)] + \ln Y_{j+1}.$$

while

$$X(\tau_1) = (\mu - \frac{1}{2}\sigma^2)\tau_1 + \sigma W(\tau_1) + \ln Y_1,$$

where  $W(t)$ ,  $\{(\tau_i, Y_i)\}$  are mutually independent. Thus, once  $\{(\tau_i, Y_i)\}$  are generated,  $X(\tau_j)$  can be generated recursively. Note that in this case,  $\tau_i$  and  $Y_i$  need not be independent.

## Gamma processes

For the jump-diffusion with  $J(t)$  a compound Poisson process,  $X(t) = \ln S(t)$  is a process with independent increments. Processes with independent increments are called Lévy processes and have wide applications. Such processes  $X(t)$ , for example, can be used to model the asset prices

$$S(t) = S(0)e^{X(t)}.$$

Under some technical conditions, a Lévy process can be represented as the sum of a diffusion and a “pure-jump” process independent of the diffusion. A pure-jump process is one whose increments are made purely of jumps. If the number of jumps in every finite interval is almost surely finite, then the pure-jump process is a compound Poisson process. On the other hand, there are many pure-jump process that has infinitely many jumps between every  $a < b$ .

A Gamma process  $X(t)$  is a Lévy process such that  $X(0) = 0$  and for any  $s < t$ ,

$$X(t) - X(s) \sim \text{Gamma}(a(t-s), \beta)$$

where  $(a, \beta)$  are the parameters of the process. It is seen that the process has stationary increments. From its definition, it is unclear whether such a process actually exists, and whether it has infinitely many jumps between any  $a < b$ . This becomes clear once the structure of the Gamma distributions are understood.

Clearly, a Gamma process is nondecreasing. One can use instead

$$X(t) = U(t) - D(t)$$

to model a process that has ups and downs, where  $U(t)$  and  $D(t)$  are two independent Gamma processes.

It is interesting that if  $U(1)$  and  $D(1)$  follow the same distribution, then  $X(t)$  can be also represented as

$$X(t) = W(G(t)),$$

where  $W$  is a standard BM and  $G$  a Gamma process independent of  $W$ . Thus,  $X$  can be viewed as the result of applying a random time-change to  $W$ . Given  $G(t)$ , the variance of  $X(t)$  is  $G(t)$ . Therefore,  $X$  is also named variance gamma process.

Like BM, the Gamma process can be simulated using random walk construction and bridge construction. The former is straightforward. The bridge construction is based on the following fact. For  $0 < t_1 < \dots < t_n$ , conditioning on the values of  $X(t_i)$ ,  $i = 1, \dots, n$ , the Gamma processes on the intervals  $(0, t_1)$ ,  $(t_1, t_2)$ ,  $\dots$ ,  $(t_{n-1}, t_n)$ , and  $(t_n, \infty)$  are independent of each other. Assume that  $X(t) \sim \text{Gamma}(at, \beta)$ . Then conditioning on  $X(t_{i-1}) = x$ , and  $X(t_i) = y$ , for any  $t_{i-1} < s < t_i$ ,

$$\frac{X(s) - x}{y - x} \sim \text{Beta}(a(s - t_1), a(t_2 - s)).$$

In other words, one can first draw  $Z \sim \text{Beta}(a(s - t_1), a(t_2 - s))$  and then set

$$X(s) = (1 - Z)x + Zy.$$



## Monte Carlo Integration

# Basics of Monte Carlo Integration

Monte Carlo is in particular useful for multidimensional problems that require the estimation of

$$\mu = \mathbb{E}h(\mathbf{X}),$$

where  $\mathbf{X}$  is a random vector and  $h$  a function.

The form of a MC method is simple. Basically, it will sample  $\mathbf{X}_1, \dots, \mathbf{X}_n$  iid  $\sim \mathbf{X}$  and approximate  $\mu$  by

$$\hat{\mu}_n = \frac{1}{n} \sum_{i=1}^n h(\mathbf{X}_i).$$

- ①  $\hat{\mu}_n$  is *consistent*: by the Strong Law of Large Numbers, with probability 1,

$$\hat{\mu}_n \rightarrow \mathbb{E}h(\mathbf{X}), \quad \text{as } n \rightarrow \infty.$$

- ②  $\hat{\mu}_n$  is *unbiased*:

$$\mathbb{E}(\hat{\mu}_n) = \mathbb{E}h(\mathbf{X}).$$

- ③ Variance

$$\text{Var}(\hat{\mu}_n) = \frac{\text{Var}[h(\mathbf{X})]}{n} = \frac{\sigma^2}{n}$$

and  $\sigma^2$  can be estimated by

$$\widehat{\text{Var}}(\hat{\mu}_n) = \frac{1}{n-1} \sum_{i=1}^n [h(\mathbf{X}_i) - \hat{\mu}_n]^2.$$

# MC for Evaluating Integrals

MC essentially is a method to numerically integrate functions. It therefore can be used in problems that have no obvious connection to random variables. A typical case is the evaluation of a “usual” integral,  $\int H(x) dx$ , where  $x$  may be of multidimensional. A general idea is to factorize  $H(x) = f(x)h(x)$ , where  $f(x)$  is a pdf, and then approximate the integral by

$$\int H(x) dx \approx \frac{1}{n} \sum_{i=1}^n h(X_i), \quad X_1, \dots, X_n \text{ iid } \sim f(x).$$

For example, to compute

$$I = \int_{-\infty}^{\infty} \ln |x| e^{-(x+1)^2/8} dx,$$

set  $h(x) = \sqrt{8\pi} \ln |x|$  and  $f(x) = e^{-(x+1)^2/8} / \sqrt{8\pi}$ . Since  $f(x)$  is the density of  $N(-1, 4)$ ,

$$I \approx \frac{\sqrt{8\pi}}{n} \sum_{i=1}^n \ln |X_i|, \quad X_i \sim N(-1, 4).$$

Generally speaking, the choice of the factorization  $H(x) = f(x)h(x)$  is critical to how well the MC approximation works.

# MC for European Call Option

Let

$S(t)$  = price of a stock at time  $t$ .

A European call option grants its holder the right to buy the stock at a fixed price (“strike price”)  $K$  at time  $T$  from now.

- If  $S(T) > K$ , the holder exercises the option to buy the stock at price  $K$  and immediately sell it to realize a profit of  $S(T) - K$ .
- If  $S(T) \leq K$ , the holder lets the option expire worthless.

The payoff at time  $T$  is

$$[S(T) - K]^+ = \max\{S(T) - K, 0\}.$$

If the continuously compounded interest rate is  $r$ , then the expected present value of the option is

$$P = \mathbb{E}\{e^{-rT}[S(T) - K]^+\}.$$

Under the Black-Scholes model for  $S$ ,

$$\frac{dS(t)}{S(t)} = r dt + \sigma dW(t),$$

where  $W$  is a standard Brownian motion, one has

$$S(T) = S(0)e^{(r-\frac{1}{2}\sigma^2)T+\sigma W(T)}.$$

Since

$$W(T) = \sqrt{T}Z, \quad \text{with } Z \sim N(0, 1),$$

$P$  can be evaluated explicitly. Nevertheless, as an illustration of the Monte Carlo method,  $P$  can be approximated by

$$\frac{e^{-rT}}{n} \sum_{i=1}^n \left[ S(0)e^{(r-\frac{1}{2}\sigma^2)T+\sigma\sqrt{T}Z_i} - K \right]^+$$

with  $Z_1, \dots, Z_n$  iid  $\sim N(0, 1)$ .

## Path-dependent MC

To value more complicated derivative securities with more complicated models of  $S(t)$ , one often has to simulate  $S(t)$  over multiple intermediate dates.

- 1 The payoff of a derivative security may depend explicitly on the values of  $S(t)$  at multiple dates.
- 2 The transitions of  $S(t)$  may not be known exactly and hence  $[0, T]$  has to be divided into smaller subintervals to obtain a reasonable approximation to the distribution of  $S(T)$ .

*Asian options* are path-dependent options. Their payoffs depend on the average level of  $S(t)$ , for example, the payoff  $(\bar{S} - K)^+$  with

$$\bar{S} = \frac{1}{m} \sum_{j=1}^m S(t_j)$$

for some fixed set of dates  $0 = t_0 < t_1 < \dots < t_m = T$ .



To estimate  $\mathbb{E}[(\bar{S} - K)^+]$ , a simple MC procedure is as follows.

- 1 For  $i = 1, \dots, N$ , simulate  $S(t_1^{(i)}), \dots, S(t_m^{(i)})$  and set

$$\bar{S}_i = \frac{1}{m} \sum_{j=1}^m S(t_j^{(i)})$$

- 2 The average of  $(\bar{S}_1 - K)^+, \dots, (\bar{S}_N - K)^+$  is taken as an estimate of  $\mathbb{E}[(\bar{S} - K)^+]$ .

Similarly, MC procedure can be used to price other path dependent options.

- Barrier options. A typical example is one that gets terminated worthless if  $S(t)$  crosses a level specified beforehand. For instance, a *down-and-out call* with barrier  $b$ , strike  $K$  and expiration  $T$  has payoff

$$1_{\{\text{all } S(t_i) \geq b\}} (S(T) - K)^+.$$

- Lookback options. In calls expiring at  $T = t_n$ , the strike price is  $\min S(t_i)$ , so the payoff is then  $S(t_n) - \min S(t_i)$ . In lookback puts, the strike price is  $\max S(t_i)$ , so the payoff is  $\max S(t_i) - S(t_n)$ .

## Evaluating Ratio of Normalizing Constants

Suppose a target density  $p$  is known up to a hard-to-know factor  $Z$ , i.e.,  $p(\mathbf{x}) = q(\mathbf{x})/Z$  with

$$Z = \int q(\mathbf{x}) d\mathbf{x}.$$

We have seen procedures that are designed to draw from  $p$  without having to evaluate  $Z$ . However, in many cases,  $Z$  has to be taken into consideration one way or another other. For example, in likelihood ratio test, one is interested in

$$\frac{p_2(\mathbf{x})}{p_1(\mathbf{x})}$$

where  $p_i(\mathbf{x})$  are likelihoods of data  $\mathbf{x}$  under two competing hypotheses. Again, in many cases, it is known that  $p_i(\mathbf{x}) \propto q_i(\mathbf{x})$ . Because

$$\frac{p_2(\mathbf{x})}{p_1(\mathbf{x})} = \frac{q_2(\mathbf{x})/Z_2}{q_1(\mathbf{x})/Z_1} = \frac{q_2(\mathbf{x})}{q_1(\mathbf{x})} \frac{Z_1}{Z_2}, \quad \text{where} \quad Z_i = \int q_i(\mathbf{x}) d\mathbf{x}.$$

it is necessary to evaluate the ratio  $Z_2/Z_1$ .

Ideally, a MC integration procedure will *not* evaluate  $Z_i$  separately. A relatively straightforward solution is as follows. We know that using the rejection method, one can sample  $\mathbf{X}_i$  iid from  $p_1$ . Let  $\mathbf{X}_1, \dots, \mathbf{X}_n$  be an iid sample from  $p_1$ . Then a MC estimate of  $(Z_2/Z_1)$  is

$$\widehat{(Z_2/Z_1)} = \frac{1}{n} \sum_{i=1}^n \frac{q_2(\mathbf{X}_i)}{q_1(\mathbf{X}_i)}.$$

This follows from

$$\mathbb{E} \left[ \frac{q_2(\mathbf{X})}{q_1(\mathbf{X})} \right] = \int \frac{q_2(\mathbf{x})}{q_1(\mathbf{x})} p_1(\mathbf{x}) d\mathbf{x} = \frac{Z_2}{Z_1} \int p_2(\mathbf{x}) d\mathbf{x} = \frac{Z_2}{Z_1}. \quad \square$$

# Bayesian Missing Data Problem

In a missing data problem, suppose  $\mathbf{Y} \sim p_{\mathbf{Y}}(\mathbf{y} | \boldsymbol{\theta})$  but only  $\mathbf{X} = M(\mathbf{Y})$  is observed, where  $M$  is a many-to-fewer mapping.

In the frequentist approach to the problem, in order to estimate  $\boldsymbol{\theta}$  by the MLE method, at least, one has to evaluate

$$L(\boldsymbol{\theta} | \mathbf{x}) = p_{\mathbf{X}}(\mathbf{x} | \boldsymbol{\theta}) = \int_{M(\mathbf{y})=\mathbf{x}} p_{\mathbf{Y}}(\mathbf{y} | \boldsymbol{\theta}) d\mathbf{y}.$$

The integral is often difficult to compute analytically.

Using MC integration, one designs a suitable function  $h$  and a distribution  $P$ , then samples  $\mathbf{Y}_1, \dots, \mathbf{Y}_n$  from  $P$  to get

$$L(\boldsymbol{\theta} | \mathbf{x}) \approx \frac{1}{n} \sum_{i=1}^n h(\mathbf{Y}_i).$$

In the Bayesian approach to the problem, let  $\theta$  be assigned with a prior distribution  $\pi(\theta)$ . One is interested in the posterior distribution of  $\theta$  when the data is observed. If the full data  $\mathbf{Y} = \mathbf{y}$  were observed, then we can evaluate the posterior distribution of  $\theta$  given  $\mathbf{y}$ . By Bayes formula,

$$f(\theta | \mathbf{y}) \propto p_{\mathbf{Y}}(\mathbf{y} | \theta)\pi(\theta).$$

If only  $\mathbf{X} = \mathbf{x}$  is observed, then the posterior distribution is

$$f(\theta | \mathbf{x}) = \int f(\theta | \mathbf{y})p(\mathbf{y} | M(\mathbf{y}) = \mathbf{x}) d\mathbf{y}.$$

Provided that  $f(\theta | \mathbf{y})$  has a closed form, a MC integration method may sample

$$\mathbf{Y}_1, \dots, \mathbf{Y}_n \sim p(\mathbf{y} | M(\mathbf{y}) = \mathbf{x})$$

and approximate

$$f(\theta | \mathbf{x}) \approx \frac{1}{n} \sum_{i=1}^n f(\theta | \mathbf{Y}_i).$$

□

# Importance Sampling, A First Look

To evaluate  $\mu = \mathbb{E}h(\mathbf{X})$ , it often works better *not* to directly draw from the distribution of  $\mathbf{X}$ , but instead to draw from a different, auxiliary distribution. Assume  $\mathbf{X} \sim p$ . Let  $q$  be another distribution, such that  $q(\mathbf{x}) > 0$  wherever  $p(\mathbf{x}) > 0$ . Then we can approximate  $\mu$  by

$$\mu \approx \frac{1}{n} \sum_{i=1}^n \frac{h(\mathbf{Y}_i)p(\mathbf{Y}_i)}{q(\mathbf{Y}_i)},$$

with  $\mathbf{Y}_1, \dots, \mathbf{Y}_n \sim q$ , because

$$\begin{aligned} \mathbb{E} \left[ \frac{h(\mathbf{Y})p(\mathbf{Y})}{q(\mathbf{Y})} \right] &= \int \frac{h(\mathbf{y})p(\mathbf{y})}{q(\mathbf{y})} q(\mathbf{y}) \, d\mathbf{y} \\ &= \int h(\mathbf{y})p(\mathbf{y}) \, d\mathbf{y} = \mathbb{E}[h(\mathbf{X})]. \end{aligned}$$

This method is called “*importance sampling*”, which is an important method to reduce the variance of MC estimate.

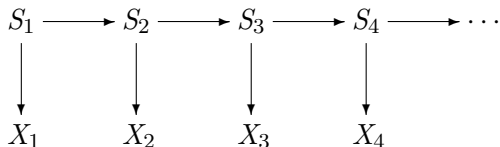
# State-Space Model

The state-space model is a dynamic system that consists of two parts:

$$\text{(state equation): } S_t \sim q_t(s_t | S_{t-1}, \theta)$$

$$\text{(observation equation): } X_t \sim f_t(x_t | S_t, \phi)$$

- ①  $S_t$  and  $X_t$  can be multi-dimensional.
- ② Only  $X_1, X_2, \dots$  are observed,  $S_1, S_2, \dots$  are hidden.
- ③  $\theta$  and  $\phi$  are system parameters.





Non-linear filtering aims to estimate, for  $t \geq 1$ , the value of  $S_t$  given the up-to-date history of the observations

$$X_1 = x_1, \dots, X_t = x_t.$$

For example, if the system parameters are known, then the least-square-error “on-line” estimate of  $S_t$  is the Bayes solution

$$\hat{s}_t = \mathbb{E}[S_t | x_1, \dots, x_t] = \int s_t p(s_t | x_1, \dots, x_t) ds_t.$$

Depending on questions being asked, median, mode, etc. of the conditional distribution of  $S_t$  given  $x_1, \dots, x_t$  can also be used as estimators. The question is how to *compute* a given estimator.

- 1 Linear state-space model: both  $f_t$  and  $q_t$  are linear Gaussian conditional distributions, i.e.,

$$S_t = A_t S_{t-1} + \varepsilon_t, \quad X_t = B_t S_t + \delta_t$$

where  $A_t$ ,  $B_t$  are matrices,  $\varepsilon_t$  and  $\delta_t$  are independent Gaussian random vectors with mean 0.

In this case, the estimate can be obtained analytically, and the resulting algorithm is the famous Kalman filter.

- 2 Hidden Markov model (HMM):  $S_t$  can only take a finite number of values, or 'states'.

Dynamic programming method can be used to evaluate the integrals.

Other than these two cases, the evaluation of the integrals is difficult and in most cases has to resort to MC. A popular algorithm based on MC integration is the so-called "particle filtering".

# Partical Filtering

Henceforth, write

$$x_{1:t} = (x_1, \dots, x_t), \quad s_{1:t} = (s_1, \dots, s_t), \quad \text{etc.}$$

Particle filtering at each time  $t$  keeps a finite set  $s_t^{(1)}, \dots, s_t^{(m)}$  to approximately follow the analytically intractable  $p(s_t | x_{1:t})$ . Each  $s_t^{(j)}$  is referred to as a “particle”. Then mean, median, etc. can be estimated using the particles, for example,

$$\mathbb{E}[s_t | x_{1:t}] \approx \text{Mean of } s_t^{(1)}, \dots, s_t^{(m)},$$

$$\text{Median}[s_t | x_{1:t}] \approx \text{Median of } s_t^{(1)}, \dots, s_t^{(m)}.$$

For  $t = 1$ ,

$$p(s_1 | x_1) \propto f_1(x_1 | s_1) q_1(s_1)$$

is typically tractable. So one can set

$$s_1^{(1)}, \dots, s_1^{(m)} \text{ i.i.d. } \sim p(s_1 | x_1)$$

The question is how to update the set of particles to approximate  $p(s_{t+1} | x_{1:t+1})$  when  $x_{t+1}$  comes in. First,

$$p(s_{t+1} | x_{1:t+1}) = \int p(s_t, s_{t+1} | x_{1:t}, x_{t+1}) ds_t.$$

Apply Bayesian formula to  $s_t$ ,  $s_{t+1}$  and  $x_{t+1}$ , while treating  $x_{1:t}$  as a fixed parameter. Then

$$p(s_t, s_{t+1} | x_{1:t}, x_{t+1}) \propto f(x_{t+1} | s_t, s_{t+1}, x_{1:t})p(s_{t+1} | s_t, x_{1:t})p(s_t | x_{1:t}).$$

By the state equation and observation equation, the right hand side is

$$f(x_{t+1} | s_{t+1})p(s_{t+1} | s_t)p(s_t | x_{1:t}).$$

Then

$$p(s_{t+1} | x_{1:t+1}) \propto f_{t+1}(x_{t+1} | s_{t+1}) \int q_{t+1}(s_{t+1} | s_t)p(s_t | x_{1:t}) ds_t$$

If we write

$$p(s_{t+1} | x_{1:t}) = \int q_{t+1}(s_{t+1} | s_t) p(s_t | x_{1:t}) ds_t, \quad (\dagger)$$

then

$$p(s_{t+1} | x_{1:t+1}) \propto f_{t+1}(x_{t+1} | s_{t+1}) p(s_{t+1} | x_{1:t}).$$

Therefore, if we can sample from  $p(s_{t+1} | x_{1:t})$ , we can sample from  $p(s_{t+1} | x_{1:t+1})$  using rejection sampling or SIR. Now if  $p(s_t | x_{1:t})$  can be sampled, then  $p(s_{t+1} | x_{1:t})$  can be sampled. This leads to

- 1) sample  $s_{t+1}^*$  from  $p(s_{t+1} | x_{1:t})$  as follows
  - (a) sample  $s_t$  from  $p(s_t | x_{1:t})$ ;
  - (b) given  $s_t$ , sample  $s_{t+1}^*$  from  $q_{t+1}(s_{t+1} | s_t)$ ;
- 2) fixing a constant  $a$  such that  $f_{t+1}(x_{t+1} | s_{t+1}) \leq a$  for all possible  $s_{t+1}$ , accept  $s_{t+1}^*$  with probability  $f_{t+1}(x_{t+1} | s_{t+1}^*)/a$ , otherwise go back to step 1.

The hope is that at each step  $t$ ,  $s_t^{(1)}, \dots, s_t^{(m)}$  approximately follow  $p(s_t | x_{1:t})$ . Then one can repeat the above procedure, except that in step 1),  $s_t$  is sampled uniformly from  $s_t^{(1)}, \dots, s_t^{(m)}$ . This will generate a new set of particles  $s_{t+1}^{(1)}, \dots, s_{t+1}^{(m)}$  that hopefully approximately follow  $p(s_{t+1} | x_{1:t+1})$ . The procedure can then continue.

Alternatively, SIR can be used:

- 1 For each  $j = 1, \dots, m$ , sample  $s_{t+1}^{(j*)}$  from  $q_{t+1}(s_{t+1} | s_t^{(j)})$ .
- 2 Generate  $s_{t+1}^{(1)}, \dots, s_{t+1}^{(m)}$  by resampling from  $s_{t+1}^{(1*)}, \dots, s_{t+1}^{(m*)}$  with replacement, with probabilities proportional to  $f_{t+1}(x_{t+1} | s_{t+1}^{(j*)})$ . □

The same procedure can be used for updating prediction. Here one needs to sample  $p(s_{t+2} | x_{1:t+1})$ . This can be done by

$$p(s_{t+2} | x_{1:t+1}) \propto \int p(s_{t+2} | s_{t+1}) p(s_{t+1} | x_{1:t+1}) ds_{t+1}.$$

The formula has already occurred in ( $\dagger$ ).

If after estimating  $S_t$  or predicting  $S_{t+1}$  based on  $X_{1:t}$ , an action is done to modify the parameters of the state equation, then the system has feedback. In this case, the above procedure still applies.

# Computing time for unbiased estimator

Three important considerations

- Computing time
- Bias
- Variance.

Suppose

$$\hat{c}_n = \frac{1}{n} \sum_{i=1}^n C_i$$

is an unbiased estimator of  $c$ , with

$$C_i \text{ iid, } \mathbb{E}(C_i) = c, \quad \text{Var}(C_i) = \sigma^2 < \infty.$$

By CLT,  $\sqrt{n}(\hat{c}_n - c) \Rightarrow N(0, \sigma^2)$ , so roughly speaking, with  $n$  replications, the magnitude of error in the estimator  $\hat{c}_n$  is about  $\sigma/\sqrt{n}$ .



Suppose  $s$  is the available time budget to produce an estimate. Since generating each  $C_i$  takes some time  $\tau_i$ , the number of  $C_i$  that can be generated is finite. The issue is how  $\tau_i$  affects the accuracy of  $\hat{c}_n$ .

Case 1:  $\tau_i \equiv \tau$  nonrandom. Since  $n = \lfloor s/\tau \rfloor$ , for  $s \gg 0$ , the magnitude of error in  $\hat{c}_n$  is about  $(\sigma/\sqrt{s})\sqrt{\tau}$ . So, given a fixed large  $s$ , the magnitude of error is proportional to  $\sqrt{\tau}$ .

Case 2:  $\tau_i$  is random. In many cases,  $\tau_i$  can vary substantially across replications. With  $\tau_i$  being random,  $n$  is random as well:

$$n = N(s) = \max \{n \geq 0 : \tau_1 + \cdots + \tau_n \leq s\}.$$

Since  $\tau_i$  are iid  $\sim \tau$ , if  $\mathbb{E}\tau < \infty$ , then for  $s \gg 0$ ,

$$N(s) \approx s/\mathbb{E}\tau$$

and hence the magnitude of error  $\approx (\sigma/\sqrt{s})\sqrt{\mathbb{E}\tau}$ .

# Bias

Some MC are biased for all finite sample sizes but become asymptotically unbiased as  $n \rightarrow \infty$ . For these estimators, the bias is negligible.

Example: For random  $\tau_i$ , in general  $\mathbb{E}[\hat{C}_{N(s)}] \neq c$ , where  $s$  is the available running time. Since

$$\sqrt{s}(\hat{C}_{N(s)} - C) \Rightarrow N(0, \sigma^2 \mathbb{E}\tau)$$

as  $s \rightarrow \infty$ , the bias decreases in the same rate as  $1/\sqrt{s}$ .

Example: Let  $(X_i, Y_i)$  be iid  $\sim (X, Y)$ . Then

$$\frac{\bar{X}_n}{\bar{Y}_n} \text{ is a biased estimator of } \frac{\mathbb{E}(X)}{\mathbb{E}(Y)}$$

for all  $n$  because

$$\mathbb{E} \left[ \frac{\bar{X}_n}{\bar{Y}_n} \right] \neq \frac{\mathbb{E}(\bar{X}_n)}{\mathbb{E}(\bar{Y}_n)} = \frac{\mathbb{E}(X)}{\mathbb{E}(Y)}.$$

However, actually, as  $n \rightarrow \infty$

$$\sqrt{n} \left[ \frac{\bar{X}_n}{\bar{Y}_n} - \frac{\mathbb{E}(X)}{\mathbb{E}(Y)} \right]$$

is asymptotically normal.

Example: Let  $f_a(x)$  be a set of functions parametrized by  $a \in A$ . In various optimization problems, it is of interest to find

$$a_0 = \arg \min_{a \in A} \mathbb{E}[f_a(X)].$$

A simple idea is to draw  $X_1, \dots, X_n$  iid, and estimate  $a_0$  by

$$\hat{a}_0 = \arg \min_{a \in A} \frac{1}{n} \sum_{i=1}^n f_a(X_i).$$

This has to be dealt with carefully because of the bias

$$\mathbb{E} \left[ \min_{a \in A} \frac{1}{n} \sum_{i=1}^n f_a(X_i) \right] \leq \min_{a \in A} \mathbb{E}[f_a(X)].$$

As long as  $A$  is finite, the bias vanishes as  $n \rightarrow \infty$ . However, if  $f_a(x)$  is difficult to evaluate, the potential bias may be difficult to remove practically as it is unpractical to evaluate  $f_a(x)$  for a large number of  $X_i$ .

Often, however, the bias does not vanish as  $n \rightarrow \infty$ .

Example: In simulating

$$S(t_j), \quad t_j = \frac{jT}{m}, \quad j = 0, 1, \dots, m,$$

under the Black-Sholes model, if one uses the naive approximation

$$S(t_{j+1}) = S(t_j) + rS(t_j)h + \sigma S(t_j)\sqrt{h}Z_{j+1},$$

the expected discounted payoff of an Asian call option will differ from the exact value. The bias typically vanishes as

$$h = T/m \rightarrow 0.$$

However, decreasing  $h$  leads to higher computational burden.

As the next example shows, if two unbiased estimators are combined to estimate another quantity, the composite estimator may be biased.

Example: Suppose  $A$  is a call option expiring at  $T$  with strike price  $K_A$ , and  $B$  is a (compound) option expiring at  $t$  with strike price  $K_B$  to buy  $A$ :

$$S(0) \longrightarrow \underset{B}{S(t)} \longrightarrow \underset{A}{S(T)}.$$

Let  $P_A(x)$  denote the expected discounted payoff of  $A$  given that  $S(t) = x$ :

$$P_A(x) = \mathbb{E}\{e^{-r(T-t)}[S(T) - K_A]^+ \mid S(t) = x\}.$$

The expected discounted payoff of  $B$  is then

$$P_B = \mathbb{E}\{e^{-rt}[P_A(S(t)) - K_B]^+\}.$$

To estimate  $P_B$ , one can draw  $m$  replicates  $s_i$  of  $S(t)$  to get an estimate

$$\frac{e^{-rt}}{n} \sum_{i=1}^n [P_A(s_i) - K_B]^+.$$

If, for each value  $s$  of  $S(t)$ ,  $P_A(s)$  can be evaluated exactly using a closed-form formula, then the estimate is unbiased.

However, if  $P_A(s)$  does not have a closed-form formula, then one has to draw  $x_j(s)$ ,  $j = 1, \dots, k$ , from the *conditional* distribution of  $S(T)$  given  $S(t) = s$ , to get an unbiased estimate

$$\hat{P}_A(s) = \frac{e^{-r(T-t)}}{k} \sum_{j=1}^k [x_j(s) - K_A]^+.$$

The estimator of  $P_B$  is a compound of the two unbiased estimators

$$\hat{P}_B = \frac{e^{-rt}}{n} \sum_{i=1}^n [\hat{P}_A(s_i) - K_B]^+,$$

often referred to as a “plug-in” estimator.

It turns out that  $\hat{P}_B$  has a positive bias,

$$\mathbb{E}(\hat{P}_B) \geq \mathbb{E}(P_B),$$

with “=” only under special circumstances. This means if someone uses  $\hat{P}_B$  to price the compound option, he will overprice it. The bias can only be reduced by increasing  $k$  but not  $n$ .

To see how the positive bias comes, recall Jensen's inequality: given random variables  $X$  and  $Y$  and convex function  $h$ ,

$$\mathbb{E}[h(X) | Y] \geq h[\mathbb{E}(X | Y)]$$

and hence

$$\mathbb{E}[h(X)] = \mathbb{E} \{ \mathbb{E}[h(X) | Y] \} \geq \mathbb{E} \{ h[\mathbb{E}(X | Y)] \}.$$



Let  $X = \hat{P}_A(S(t))$ ,  $Y = S(t)$ ,  $h(x) = e^{-rt}(x - K_B)^+$ . Then

- ①  $h(x)$  is convex,

$$P_B = \mathbb{E}[h(P_A(S(t)))].$$

and

$$\hat{P}_B = \frac{1}{n} \sum_{i=1}^n h(\hat{P}_A(s_i)) = \frac{1}{n} \sum_{i=1}^n h(X_i),$$

where  $s_1, \dots, s_n$  iid  $\sim S(t)$  and  $X_i = \hat{P}_A(s_i)$ .

- ② For each  $s_i$ ,  $\hat{P}_A(s_i)$  is an unbiased estimator of  $P_A(s)$ :

$$\mathbb{E}[X_i] = P_A(s_i).$$

Then

$$\mathbb{E}[\hat{P}_B] = \mathbb{E}[h(X)] \geq \mathbb{E}\{h(\mathbb{E}(X | Y))\} = \mathbb{E}[h(P_A(S(t)))] = P_B. \quad \square$$

# Mean squared error

A standard measure of performance is mean squared error (MSE). If  $\hat{c}$  is an estimator of  $C$ , then

$$\begin{aligned}\text{MSE}(\hat{c}) &= \mathbb{E}[(\hat{c} - C)^2] \\ &= [\mathbb{E}(\hat{c}) - C]^2 + \mathbb{E}[(\hat{c} - \mathbb{E}(\hat{c}))^2] \\ &= \text{Bias}^2(\hat{c}) + \text{Var}(\hat{c}).\end{aligned}$$