

2. Properties of the finite-dimensional Dirichlet distribution.

The Dirichlet distribution is one of the possible multivariate extension of the Beta. We start considering:

$$U_1, \dots, U_k \text{ s.t. } U_j \sim \text{Gamma}(\alpha_j, \beta) \quad j=1, \dots, k$$

$$\text{We define: } X_j = \frac{U_j}{U_1 + \dots + U_k}$$

We obtain: $(X_1, \dots, X_k) \sim \text{Dirichlet}(\alpha_1, \dots, \alpha_k)$.

We impose $\sum_{j=1}^k U_j = 1$ and so $(X_1, \dots, X_k) \in \Delta_{k-1} = \{\underline{x} \in \mathbb{R}^k : \sum_{j=1}^k x_j = 1, x_j \in [0, 1] \forall j\}$

The p.d.f. of the Dirichlet distribution is:

$$f(x_1, \dots, x_k; \alpha_1, \dots, \alpha_k) = \frac{1}{B(\alpha)} \prod_{j=1}^k x_j^{\alpha_j - 1} \mathbf{1}_{\Delta_{k-1}}(\underline{x})$$

$$\text{where } B(\alpha) = \frac{\Gamma(\alpha_1) \dots \Gamma(\alpha_k)}{\Gamma(\alpha_1 + \dots + \alpha_k)}$$

Notice that, because of Δ_{k-1} : $f(\underline{x}; \alpha) = \frac{1}{B(\alpha)} \left[\prod_{j=1}^{k-1} x_j^{\alpha_j - 1} \right] \left(1 - \sum_{j=1}^{k-1} x_j \right)^{\alpha_k - 1} \mathbf{1}_{\Delta_{k-1}}(\underline{x})$

$$\text{For simplicity: } \alpha_0 := \sum_{j=1}^k \alpha_j$$

Properties:

$$1. \mathbb{E}[X_j] = \frac{\alpha_j}{\alpha_0} \quad j=1, \dots, k$$

$$2. \text{Var}(X_j) = \frac{\alpha_j (\alpha_0 - \alpha_j)}{\alpha_0^2 (\alpha_0 + 1)} \quad j=1, \dots, k$$

$$3. \text{Cov}(X_i, X_j) = -\frac{\alpha_i \alpha_j}{\alpha_0^2 (\alpha_0 + 1)} \quad i \neq j$$

$$4. \text{Marginal distribution: } X_j \sim \text{Beta}(\alpha_j, \alpha_0 - \alpha_j) \quad j=1, \dots, k$$

5. Aggregation property:

$$(X_1, \dots, X_k) \sim \text{Dirichlet}(\alpha_1, \dots, \alpha_k)$$

$$(Z_1, \dots, Z_m) := (\sum_{j \in C_1} X_j, \dots, \sum_{j \in C_m} X_j) \quad \text{s.t.} \quad \begin{cases} C_j \subseteq \{1, \dots, n\} & \forall j \\ C_i \cap C_j = \emptyset & \forall i \neq j \\ \bigcup_{j=1}^m C_j = \{1, \dots, n\} \end{cases}$$

$$\rightarrow (Z_1, \dots, Z_m) \sim \text{Dirichlet}(\alpha_1^*, \dots, \alpha_m^*)$$

$$\alpha_j^* = \sum_{j \in C_j} \alpha_j$$

6. If $(X_1, \dots, X_k) \sim \text{Dirichlet}(\alpha_1, \dots, \alpha_k)$ then:

$$\begin{cases} X_1 \sim \text{Beta}(\alpha_1, \alpha_0 - \alpha_1) \\ X_2 \perp\!\!\!\perp \left(\frac{X_2}{X_1}, \dots, \frac{X_k}{X_1} \right) \\ \left(\frac{X_2}{X_1}, \dots, \frac{X_k}{X_1} \right) \sim \text{Dirichlet}(\alpha_2, \dots, \alpha_k) \end{cases} \quad \leftarrow \text{we can take out one dimension}$$

3. Bayes' Theorem for dominated models

Consider a random sample $X_n|\theta = (X_1, \dots, X_n)|\theta \sim P_\theta$, $\theta \in \Theta \subset \mathbb{R}^n$. Suppose that the model is dominated (\exists a density (the R-N density) : $f(\underline{x}|\theta) = \text{density of the joint distribution of } X_n \text{ w.r.t. some } \sigma\text{-finite measure on } \mathbb{R}^n (\therefore = \lambda^{(n)})$). Suppose $\theta \sim \pi$ (prior for θ).

\Rightarrow The posterior distribution for θ given $X_n = \underline{x}$ can be expressed as:

$$P(\theta \in B | X_n = \underline{x}) = \frac{\int_B f(\underline{x}|\theta) \pi(d\theta)}{\int_{\Theta} f(\underline{x}|\theta) \pi(d\theta)} \quad \forall B \in \mathcal{B}(\Theta)$$

proof.

We define $\delta := \text{joint distribution of } \underline{X} \text{ and } \theta$:

$$\begin{aligned} \delta(A \times B) &= \int_B P_\theta(A) \pi(d\theta) * \int_B \int_A f(\underline{x}|\theta) \lambda^{(n)}(d\underline{x}) \pi(d\theta) \\ &\stackrel{\text{FT}}{=} \int_A \left(\int_B f(\underline{x}|\theta) \pi(d\theta) \right) \lambda^{(n)}(d\underline{x}) \end{aligned} \quad A \in \mathcal{B}(\mathbb{R}^n) \quad B \in \mathcal{B}(\Theta)$$

* is valid because P_θ is dominated \Rightarrow we can express it as integral of the density.

We want now to define the marginal:

$$\mu_n(A) = \delta(A \times \Theta) = \int_A \int_{\Theta} f(\underline{x}|\theta) \pi(d\theta) \lambda^{(n)}(d\underline{x}) := \int_A m_n(\underline{x}) \lambda^{(n)}(d\underline{x}).$$

We needed this marginal because the theorem is obtained from:

$$\underline{\chi}(\underline{X}, \theta) = \underline{\chi}(\underline{X}|\theta) \underline{\chi}(\theta) = \underline{\chi}(\theta|\underline{X}) \underline{\chi}(\underline{X}) \quad \rightarrow \text{we needed } \underline{\chi}(\underline{X}).$$

We now fix $B \in \mathcal{B}(\Theta)$ and we define $\delta(A \times B)$ varying A .

$\delta(\cdot \times B)$ is a measure on $\mathcal{B}(\mathbb{R}^n)$, moreover: $\delta(\cdot \times B) \ll \mu_n$

(i.e. $\delta(\cdot \times B)$ is absolutely continuous w.r.t. μ_n):

consider A : $\mu_n(A) = 0 \Rightarrow \delta(A \times B) \stackrel{\text{def.}}{\leq} \delta(A \times \Theta) = \mu_n(A) = 0$.

We can apply the R-N theorem:

$$\forall B \text{ fixed } \exists \text{ function } \pi(\underline{x}, B): \delta(A \times B) = \int_A \pi(\underline{x}, B) \mu_n(d\underline{x})$$

$$\begin{aligned} \rightarrow \delta(A \times B) &= \int_A \pi(\underline{x}, B) m_n(\underline{x}) \lambda^{(n)}(d\underline{x}) \\ &= \int_A \pi(\underline{x}, B) \int_{\Theta} f(\underline{x}|\theta) \pi(d\theta) \lambda^{(n)}(d\underline{x}) \end{aligned}$$

Considering that: $\delta(A \times B) = P(\underline{X} \in A, \theta \in B) = \int_A P(\theta \in B | \underline{X} = \underline{x}) \mu_n(d\underline{x})$

we have that: $\pi(\underline{x}, B) = P(\theta \in B | \underline{X} = \underline{x})$:

$$\delta(A \times B) = \int_A P(\theta \in B | \underline{X} = \underline{x}) \int_{\Theta} f(\underline{x}|\theta) \pi(d\theta) \lambda^{(n)}(d\underline{x})$$

Comparing:

$$\begin{cases} \delta(A \times B) = \int_A \int_B f(\underline{x}|\theta) \pi(d\theta) \lambda^{(n)}(d\underline{x}) \\ \delta(A \times B) = \int_A P(\theta \in B | \underline{X} = \underline{x}) \int_{\Theta} f(\underline{x}|\theta) \pi(d\theta) \lambda^{(n)}(d\underline{x}) \end{cases}$$

$$\rightarrow \int_B f(\underline{x}|\theta) \pi(d\theta) = P(\theta \in B | \underline{X} = \underline{x}) \int_{\Theta} f(\underline{x}|\theta) \pi(d\theta) \quad \text{a.s. w.r.t. } \mu_n$$

$$\rightarrow P(\theta \in B | \underline{X} = \underline{x}) = \frac{\int_B f(\underline{x}|\theta) \pi(d\theta)}{\int_{\Theta} f(\underline{x}|\theta) \pi(d\theta)} \quad \text{a.s. w.r.t. } \mu_n$$

Moreover: $\mu_n(\{\underline{x}: \int_{\Theta} f(\underline{x}|\theta) \pi(d\theta) = m_n(\underline{x}) = 0\}) = 0$

If $\theta \sim \pi(\theta)$:

$$\pi(\theta | \underline{x}) = \frac{f(\underline{x}|\theta) \pi(\theta)}{\int_{\Theta} f(\underline{x}|\theta) \pi(\theta) d\theta} \quad \Rightarrow \quad \text{The posterior distribution is proportional to the likelihood times the prior}$$

1. Joint of \underline{X}, θ (I)

2. Marginal of \underline{X}

3. Radon-Nikodym: $\delta(\cdot \times B), \mu_n(\cdot)$ (B fixed)

4. Joint of \underline{X}, θ (II)

5. (I) = (II)

4. Posterior mean as the functional minimizing the posterior quadratic loss function (multivariate case), posterior median as the functional minimizing the posterior absolute value loss function.

We suppose: $X_1, \dots, X_n | \theta \sim f(x, \theta)$, $\theta \sim \pi(\theta)$ s.t. $\theta \in \Theta \subseteq \mathbb{R}^k$.

We have a set of possible choices (actions) A : $a \in A$ is a single action (in terms of point estimation). We define a loss function: $l: (\theta, a) \mapsto l(\theta, a)$. The loss function tells us how much we lose choosing the action a when the (true) state is θ .

We have: prior expectation loss : $E_{\pi}[l(\theta, a)]$
posterior expectation loss : $E_{\pi}[l(\theta, a)|X]$

The optimal choices (by Bayes) are all the a^* s.t. :

$$a^* = \arg \min_{a \in A} E_{\pi}[l(\theta, a)|X]$$

1. Quadratic loss function (multivariate): ($\Theta \subseteq \mathbb{R}^p$)

$A = \Theta$, $l(\theta, a) = (\theta - a)^T Q (\theta - a)$ Q symmetric, positive semidef. matrix.

$$\Rightarrow \tilde{\theta} = E_{\pi}[\theta|X] = \arg \min_{\theta \in \Theta} E_{\pi}[l(\theta, a)|X]$$

proof.

$$\frac{d}{da} E_{\pi}[l(\theta, a)|X] = \frac{d}{da} \int_{\Theta} (\theta - a)^T Q (\theta - a) \pi(\theta|X) d\theta := \frac{d}{da} \int_{\Theta} F(a) \pi(\theta|X) d\theta$$

$$\text{where } F(a) = \sum_{i,j=1}^p q_{ij} (\theta_i - a_i)(\theta_j - a_j), \quad q_{ij} = [Q]_{ij}$$

$$\begin{aligned} \frac{d}{da} F(a) &= - \sum_{j=1}^p q_{ej} (\theta_j - a_j) - \sum_{i=1}^p q_{ie} (\theta_i - a_i) \\ &= - 2 \sum_{j=1}^p q_{ej} (\theta_j - a_j) = [-2Q(\theta - a)]_e \end{aligned}$$

$$\Rightarrow \frac{d}{da} E_{\pi}[l(\theta, a)|X] = \int_{\Theta} -2Q(\theta - a) \pi(\theta|X) d\theta = 0$$

$$\Rightarrow 2Q \int_{\Theta} \theta \pi(\theta|X) d\theta = 2Q \int_{\Theta} a \pi(\theta|X) d\theta \Rightarrow a$$

$$\Rightarrow a = \tilde{\theta} = \int_{\Theta} \theta \pi(\theta|X) d\theta = E_{\pi}[\theta|X]$$

2. Absolute value loss function:

$$l(\theta, a) = |\theta - a| = \begin{cases} \theta - a & \theta > a \\ a - \theta & \theta \leq a \end{cases}$$

$\Rightarrow \tilde{\theta} = \text{quantile of } \pi(\theta|X) \text{ of order } \frac{1}{2} = \text{median of } \pi(\theta|X)$

proof.

$$\begin{aligned} E_{\pi}[l(\theta, a)|X] &= \int_a^{\infty} (\theta - a) \pi(\theta|X) d\theta + \int_{-\infty}^a (a - \theta) \pi(\theta|X) d\theta \\ &= \int_a^{\infty} \int_a^{\theta} dy \pi(\theta|X) d\theta + \int_{-\infty}^a \int_{\theta}^a dy \pi(\theta|X) d\theta \\ &= \int_a^{\infty} \int_y^{\infty} \pi(\theta|X) d\theta dy + \int_{-\infty}^a \int_{-\infty}^y \pi(\theta|X) d\theta dy \\ &= \int_a^{\infty} P(\theta > y|X) dy + \int_{-\infty}^a P(\theta < y|X) dy \end{aligned}$$

$$\begin{aligned} \frac{d}{da} E_{\pi}[l(\theta, a)|X] &= -P(\theta > a|X) + P(\theta < a|X) \\ &= -(1 - P(\theta \leq a|X)) + P(\theta < a|X) \\ &= 0 \end{aligned}$$

$$\Rightarrow P(\theta \leq \tilde{\theta}|X) = \frac{1}{2} \Rightarrow \tilde{\theta} = \text{median of } \pi(\theta|X)$$

5. Use the MC std error to bound (in probability) the Monte Carlo error.

All we want to know about a distribution can be achieved by simulating many draws from it. Suppose $\theta \sim \pi$ distribution (which we consider as posterior of the Bayesian model): $\theta \in \Theta \subseteq \mathbb{R}^p$. Let's consider a measurable function $h: \Theta \rightarrow \mathbb{R}$ s.t. $E_\pi[h(\theta)] < \infty$. We define $\bar{h} := \int_\Theta h(\theta) \pi(d\theta)$. We want to evaluate \bar{h} .

Strong law of Large Numbers (SLLN):

$$\text{If } \theta^{(1)}, \theta^{(2)}, \theta^{(3)}, \dots \stackrel{iid}{\sim} \pi \text{ and } \bar{h}^{(T)} := \frac{1}{T} \sum_{j=1}^T h(\theta^{(j)}) \quad \left(\begin{array}{l} h: \Theta \rightarrow \mathbb{R} \\ E_\pi[h(\theta)] < \infty \end{array} \right)$$

$$\Rightarrow \bar{h}^{(T)} \xrightarrow{T \rightarrow \infty} \bar{h} \quad \text{a.s.}$$

$$\left(\frac{1}{T} \sum_{j=1}^T h(\theta^{(j)}) \xrightarrow{T \rightarrow \infty} E_\pi[h(\theta)] = \int_\Theta h(\theta) \pi(d\theta) \right)$$

$$\text{Therefore: } \text{err}(T) := \bar{h}^{(T)} - \bar{h} \xrightarrow{T \rightarrow \infty} 0 \quad \text{a.s.}$$

The point is: how fast the error goes to zero?

Central Limit Theorem (CLT):

Let $\theta^{(1)}, \theta^{(2)}, \dots$ be an infinite sequence iid from π .

Consider $h: \Theta \rightarrow \mathbb{R}$, $\bar{h} := E_\pi[h(\theta)]$ s.t. $0 < \text{Var}(h(\theta)) < \infty$

Let's define $\sigma^2 := \text{Var}(h(\theta))$ and $\sigma^2(T) := \frac{1}{T} \sum_{j=1}^T (h(\theta^{(j)}) - \bar{h})^2$

$$\Rightarrow \begin{cases} \sqrt{T} (\bar{h}^{(T)} - \bar{h}) \xrightarrow{d} N(0, \sigma^2) \\ \sigma^2(T) \xrightarrow{T \rightarrow \infty} \sigma^2 \quad \text{a.s.} \end{cases}$$

The first conclusion is saying: $\bar{h}^{(T)} - \bar{h}$ goes to zero at $(\sqrt{T})^{-1}$

$$\Rightarrow \text{err}(T) = \bar{h}^{(T)} - \bar{h} \approx N(0, \frac{\sigma^2}{T}) \approx N(0, \frac{\sigma^2(T)}{T})$$

And so we know how to bound the error:

$$P(|\text{err}(T)| > c) = P\left(\frac{|\bar{h}^{(T)} - \bar{h}|}{\sqrt{\frac{\sigma^2(T)}{T}}} > \frac{c}{\sqrt{\frac{\sigma^2(T)}{T}}}\right) \xrightarrow{\text{CLT}} 2\left(1 - \Phi\left(\frac{c}{\sqrt{\frac{\sigma^2(T)}{T}}}\right)\right)$$

\Rightarrow The Monte Carlo standard error $\sqrt{\frac{\sigma^2(T)}{T}}$ controls in probability how large the Monte Carlo error is: the smaller the std error, the smaller the error.

6. Definition of the method of composition and its use in Bayesian statistic.

Suppose that we want to sample from $Z(X_1)$ but we're not able to.

Suppose that we can sample from $Z(X_1 | X_2)$ and $Z(X_2)$ (where X_2 is just an auxiliary random variable).

Then we recall: $Z(X_1) = \int Z(X_1, dX_2) = \int Z(X_1 | X_2) Z(dX_2)$.

And so, we can now sample from $Z(X_1)$:

\Rightarrow it's enough to sample from $Z(X_1, X_2)$: $(X_2^{(i)}, X_1^{(i)})$ MC sample from $Z(X_1, X_2)$.

How can we sample from $Z(X_1, X_2)$?

1. We sample $X_2^{(i)} \sim Z(X_2)$

2. We sample $X_1^{(i)} \sim Z(X_1 | X_2 = X_2^{(i)})$

\Rightarrow We obtain samples of the joint distribution $Z(X_1, X_2)$.

How can we obtain samples of $Z(X_1)$?

We consider only the first component of each couple.

In Bayesian Statistic we use this when we need to sample from the posterior distribution.

7. Generalized inverse distribution function sampling method:

- (i) Sampling from a truncated distribution function on a real interval $(X|_{[a,b]})$
- (ii) Sampling from a Gamma (m, b) in integer.

Inverse transform method:

We want to sample $X \sim F_X(\cdot)$ with $F_X(\cdot)$ c.d.f.: $F_X: \mathbb{R} \rightarrow [0, 1]$ monotonically increasing, right continuous, $F(-\infty) = 0$, $F(+\infty) = 1$.

Prop. If $U \sim U([0,1]) \Rightarrow X = F^{-1}(U) \sim F_X \quad (F^{-1}(u) = \inf \{x : F(x) \geq u\}$
Proof. $= \text{generalized inverse of } F$)

$$x_0 \in \mathbb{R} : \quad A_1 = \{u \in [0,1] : F^{-1}(u) \leq x_0\} \\ A_2 = \{u \in [0,1] : F(x_0) \geq u\}$$

Consider $u \in A_1 : F^{-1}(u) \leq x_0$. For right continuity: $u \leq F(x_0) \Rightarrow u \in A_2$.

Consider $u \in A_2 : u \leq F(x_0)$. By the def. of $F(\cdot)$: $F^{-1}(u) \leq x_0 \Rightarrow u \in A_1$.

$$\Rightarrow A_1 = A_2$$

$$\Rightarrow P(X \leq x_0) = P(F^{-1}(U) \leq x_0) = P(U \leq F(x_0)) = F_X(x_0)$$

construction $A_1 = A_2$ $U \sim U([0,1])$

\Rightarrow Algorithm: for $i = 1, \dots, n_{\text{rep}}$:
 (to sample $X \sim F(\cdot)$) 1. Sample $U_i \sim U([0,1])$
 2. Let $X_i = F^{-1}(U_i)$

(ii) Truncated distribution

Let $X \sim F_X(\cdot)$ absolutely continuous. Let $A := (a, b) \subseteq \mathbb{R}$. Let $Y = X|_{X \in A}$.
 We want to sample from Y .

$$F_Y(y) = F_{X|X \in A}(y) = P(X \leq y | X \in A) = \frac{P(X \leq y, X \in A)}{P(X \in A)} = \begin{cases} 0 & y \leq a \\ \frac{F_X(y) - F_X(a)}{F_X(b) - F_X(a)} & a < y \leq b \\ 1 & y > b \end{cases}$$

$$f_Y(y) = \frac{dF_Y}{dy} = \frac{1}{P(Y \in A)} f_X(y) \mathbf{1}_A(y)$$

$$\Rightarrow U = F_Y(y) = \frac{F_X(y) - F_X(a)}{F_X(b) - F_X(a)} \Rightarrow F_X(y) = [F_X(b) - F_X(a)]U + F_X(a) := U^*$$

\Rightarrow We sample $U^* \sim U([F_X(a), F_X(b)])$ and then $Y = F^{-1}(U^*)$.

(iii) Gamma distribution (Gamma (m, b))

1. Exponential:

$$X \sim \mathcal{E}(\lambda) : F(x) = 1 - e^{-\lambda x}$$

$$U = F(X) = 1 - e^{-\lambda X} \Rightarrow 1 - U = e^{-\lambda X} \Rightarrow \log(1 - U) = -\lambda X$$

$$\Rightarrow X = -\frac{1}{\lambda} \log(1 - U) \stackrel{d}{=} -\frac{1}{\lambda} \log(U)$$

We sample $U_i \sim U([0,1])$ and then $X_i = -\frac{1}{\lambda} \log(U_i)$

2. Gamma:

Remark 1. $X_1, \dots, X_m : X_i \sim \text{Gamma}(a_i, b) \Rightarrow \sum_{i=1}^m X_i \sim \text{Gamma}(\sum a_i, b)$

Remark 2. $\text{Gamma}(1, b) = \mathcal{E}(b)$

$X \sim \text{Gamma}(m, b)$:

for $i = 1, \dots, n_{\text{rep}}$:

1. Sample $Y_1, \dots, Y_m \stackrel{iid}{\sim} \text{Gamma}(1, b) = \mathcal{E}(b)$

2. Let $X_i = \sum_{j=1}^m Y_j$

8. Simulation from a univariate (Box-Müller) and multivariate Gaussian distribution.

Box-Müller method.

We want to sample from $X \sim N(0,1)$.

We start from $(X,Y) \stackrel{\text{iid}}{\sim} N(0,1)$: $f(x,y) = \frac{1}{2\pi} e^{-\frac{1}{2}(x^2+y^2)}$

We move to polar coordinates:

$$\begin{cases} \rho = x^2 + y^2 & \in [0, \infty) \\ \varphi = \tan^{-1}\left(\frac{y}{x}\right) & \in [0, 2\pi] \end{cases} \Rightarrow \begin{cases} x = \sqrt{\rho} \cos(\varphi) \\ y = \sqrt{\rho} \sin(\varphi) \end{cases}$$

$$J = \begin{bmatrix} \frac{\partial x}{\partial \rho} & \frac{\partial x}{\partial \varphi} \\ \frac{\partial y}{\partial \rho} & \frac{\partial y}{\partial \varphi} \end{bmatrix} = \begin{bmatrix} \frac{1}{2\rho} \cos(\varphi) & -\frac{1}{2\rho} \sin(\varphi) \\ -\frac{1}{2\rho} \sin(\varphi) & -\frac{1}{2\rho} \cos(\varphi) \end{bmatrix}, \quad |J| = \frac{1}{2}$$

$$\begin{aligned} \Rightarrow f(\rho, \varphi) &= f_X(\sqrt{\rho} \cos(\varphi), \sqrt{\rho} \sin(\varphi)) \cdot |J| \cdot \mathbb{1}_{(0,\infty)}(\rho) \mathbb{1}_{[0,2\pi]}(\varphi) \\ &= \frac{1}{2\pi} e^{-\frac{1}{2}\rho} \cdot \frac{1}{2} \\ &\sim U([0, \infty]) \sim \mathcal{E}\left(\frac{1}{2}\right) \end{aligned}$$

→ Algorithm: for $i = 1, \dots, n_{\text{rep}}/2$:

1. Sample $U_1, U_2 \sim U([0,1])$
2. Set $\rho = -2 \log(U_1)$, $\varphi = (2\pi) U_2$
3. Set $X = \sqrt{\rho} \cos(\varphi)$, $Y = \sqrt{\rho} \sin(\varphi)$

We obtain two independent gaussian random variables X and Y .

If we want:

- $Y \sim N(\mu, \sigma^2) \Rightarrow$ sample $X \sim N(0,1)$, set $Y = \sigma X + \mu$
- $Y \sim N_k(\mu, \Sigma)$, $\Sigma = A^T A \Rightarrow$ sample $X \sim N_k(0, I)$, set $Y = AX + \mu$

9.

Acceptance-rejection method to sample from a (univariate) target density f .

Application: sampling from $\text{Gaussian}(\alpha, \beta)$ with $\alpha > 1$ when the proposal is $\text{Gaussian}(m, b)$ with m integer. Parameters choice to maximize the efficiency.

Acceptance-rejection method.

Suppose we want to sample from $X \sim f(x)$ but we're not able to do it. We introduce an instrumental random variable $Y \sim g(y)$ s.t.:

1. $S_x \subseteq S_y \subseteq \mathbb{R}^k$
2. $f(x) \leq M g(x) \quad \forall x \quad (\exists M > 0)$

Prop. Assume $f(x)$ and $g(x)$ densities s.t. 1. and 2. holds. Then, to simulate from $f(x)$ we can simulate from $Y \sim g(y)$ and $U | Y=y \sim U([0, M \cdot g(y)])$ until $0 < U \leq f(y)$.

The proposition is saying: $\Pr(Y \leq x | U \leq \frac{f(y)}{M g(y)}) = \Pr(X \leq x)$

proof.

$$\begin{aligned} \Pr(Y \leq x | U \leq \frac{f(y)}{M g(y)}) &= \frac{\Pr(Y \leq x, U \leq \frac{f(y)}{M g(y)})}{\Pr(U \leq \frac{f(y)}{M g(y)})} = \frac{\int \Pr(Y \leq x, U \leq \frac{f(y)}{M g(y)}, Y=y) g(y) dy}{\int \Pr(U \leq \frac{f(y)}{M g(y)}, Y=y) g(y) dy} \\ &= \frac{\int \mathbb{1}_{(-\infty, x)}(y) \Pr(U \leq \frac{f(y)}{M g(y)}) g(y) dy}{\int \Pr(U \leq \frac{f(y)}{M g(y)}) g(y) dy} \\ &= \frac{\int \mathbb{1}_{(-\infty, x)}(y) \frac{f(y)}{M g(y)} g(y) dy}{\int \frac{f(y)}{M g(y)} g(y) dy} \\ &= \frac{\int \mathbb{1}_{(-\infty, x)}(y) f(y) dy}{\int f(y) dy} = \Pr(X \leq x) \quad \blacksquare \end{aligned}$$

Sampling from $f(\cdot)$ given $g(\cdot)$.

- when does it work?

$\exists M > 0$: 1., 2.

- find $M/M(\theta)$ given $g(\cdot)$:

$$\frac{f(x)}{g(x)} \leq h(x, \theta)$$

$$\frac{d}{dx} h(x, \theta) = 0 \Rightarrow x^* = \arg \max_{x \in S_x} h(x, \theta)$$

$$M(\theta) = h(\theta, x^*)$$

- acceptance probability?

$$\int \Pr\left(\frac{f(x)}{M g(x)} g(x) dx = \frac{1}{M(\theta)}\right)$$

$$\Rightarrow \theta^* = \arg \max \frac{1}{M(\theta)}$$

- (9.) The algorithm is: for $i = 1, \dots, n_{\text{rep}}$:
1. sample $Y \sim g(y)$
 2. sample $U \sim U(0,1)$
 3. if $U \leq \frac{f(y)}{Mg(y)}$ set $X_i = Y$, otherwise go to 1.

Sampling from Gamma($\alpha, 1$), $\alpha > 1$.

We want to sample from $f(x) \equiv \text{Gamma}(\alpha, 1)$. The proposal is $g(x) \stackrel{d}{=} \text{Gamma}(m, b)$ with $b = L \alpha J$. We need to find M s.t. $f \leq Mg$:

$$\frac{f}{g} = \frac{\frac{1}{\Gamma(\alpha)} x^{\alpha-1} e^{-x}}{\frac{1}{\Gamma(m)} b^m x^{m-1} e^{-bx}} = \frac{m(m)}{\Gamma(\alpha)} b^{-m} x^{\alpha-m} e^{-(1-b)x}$$

$\text{Gamma}(\alpha-m+1, 1-b)$

conditions: $\begin{cases} 1-b > 0 \\ \alpha-m+1 > 0 \end{cases} \quad (b < 1) \quad (\forall)$

For $b \in (0, 1)$ the maximum of the Gamma is in its mode: $\frac{\alpha-L\alpha J}{1-b}$

$$\Rightarrow \frac{f}{g} \leq \frac{\Gamma(L\alpha J)}{\Gamma(\alpha)} b^{-L\alpha J} \left(\frac{\alpha-L\alpha J}{1-b} \right)^{\alpha-L\alpha J} e^{-(1-b)\left(\frac{\alpha-L\alpha J}{1-b} \right)} := M(b)$$

Optimize in terms of $b = \min M(b) = \max M^{-1}(b)$. (since the probability of accepting is $1/m$ and we want to maximize it)

To maximize $M^{-1}(b)$ we consider the mode of Beta($L\alpha J + 1, L\alpha J - \alpha + 1$), which is in $\frac{L\alpha J}{\alpha}$

$$\Rightarrow b_{\text{optimal}} = \frac{L\alpha J}{\alpha} \Rightarrow M_{\text{optimal}} = M(b_{\text{optimal}})$$

10. Optimal choice for the proposal distribution of the importance sampler, as the one minimizing the variance of the estimator.

Importance sampling:

Let f_x be a density and $h(\cdot)$ a measurable function. We're interested in evaluating $E_f[h(X)] = \int h(x) f(x) dx$. If we're not able to sample from f_x (and so via Monte Carlo method), we can introduce an instrumental random variable with density $g(\cdot)$.

- We have:
- Generate: $X_1, \dots, X_m \sim g(x)$
 - $E_f[h(X)] \approx \frac{1}{m} \sum_{j=1}^m \frac{f(x_j)}{g(x_j)} h(x_j)$

$$\text{In fact: } \frac{1}{m} \sum_{j=1}^m \frac{f(x_j)}{g(x_j)} h(x_j) \xrightarrow{a.s.} E_g \left[\frac{f(x)}{g(x)} h(x) \right] = \int \frac{f(x)}{g(x)} h(x) g(x) dx = E_f[h(X)]$$

The optimal choice of $g(\cdot)$ is the one that minimizes the variance of the importance sampler estimator, i.e.:

$$g^*(x) = \frac{|h(x)| f(x)}{\int |h(t)| f(t) dt}$$

proof:

$$\text{Var}_g \left(h(x) \frac{f(x)}{g(x)} \right) = E_g \left[\left(h(x) \frac{f(x)}{g(x)} \right)^2 \right] - E_g^2 \left[h(x) \frac{f(x)}{g(x)} \right]$$

$$E_g \left[h(x) \frac{f(x)}{g(x)} \right] = \int h(x) \frac{f(x)}{g(x)} g(x) dx = \int h(x) f(x) dx \perp\!\!\!\perp g(x)$$

$\Rightarrow g(x)$ that minimizes $\text{Var}_g(\dots) = g(x)$ that minimizes $E_g[(\dots)^2]$

$$E_g \left[\left(h(x) \frac{f(x)}{g(x)} \right)^2 \right] \geq \left(E_g \left[|h(x)| \frac{f(x)}{g(x)} \right] \right)^2 = \left(\int |h(x)| f(x) dx \right)^2$$

If $g(x) = \frac{|h(x)| f(x)}{\int |h(x)| f(x) dx}$ we obtain the equality (\dagger)

and so such $g(x)$ is optimal.

$$\begin{aligned}
 \textcircled{*} \int_E P(x, A) \pi(x) dx &= \int_E \left(\int_A p(x, y) dy + r(x) \mathbb{1}_A(x) \right) \pi(x) dx \\
 &= \int_A \left(\int_E p(x, y) \pi(x) dx \right) dy + \int_A r(x) \pi(x) dx \\
 R &= \int_A \left(\int_E p(y, x) dx \right) \pi(y) dy + \int_A r(x) \pi(x) dx \\
 &\downarrow = \int_A (1 - r(y)) \pi(y) dy + \int_A r(x) \pi(x) dx = \int_A \pi(y) dy = \pi(A)
 \end{aligned}$$

11. Reversability of the Metropolis-Hastings algorithm wrt. the target distribution π .
Metropolis-Hastings:

$\pi(x)$ = target distribution wrt. a measure μ

$Q(x, dy) := q(x, y) \mu(dy)$ = transition probability proposal density

$$\alpha(x, y) := \begin{cases} \min\left(\frac{\pi(y) q(y, x)}{\pi(x) q(x, y)}, 1\right) & \text{if } \pi(x) q(x, y) > 0 \\ 1 & \text{otherwise} \end{cases}$$

= probability of accepting y (that comes from $Y \sim Q(x, \cdot)$) starting from x

$$p(x, y) := \begin{cases} q(x, y) \alpha(x, y) & x \neq y \\ 0 & x = y \end{cases}$$

$$P(x, A) := \int_A p(x, y) \mu(dy) + (r(x)) \mathbb{1}_A(x) = \text{transition probability prob. of remaining in } x$$

Reversability condition: $p(x, y) \pi(x) = p(y, x) \pi(y)$
proof.

since $\alpha(\cdot, \cdot)$ is a minimum we have 2 cases:

$$1. \frac{\pi(y) q(y, x)}{\pi(x) q(x, y)} < 1$$

$$p(x, y) \pi(x) = q(x, y) \alpha(x, y) \pi(x) = q(x, y) \frac{\pi(y) q(y, x)}{\pi(x) q(x, y)} \stackrel{\text{if } \alpha < 1}{=} \pi(y) q(y, x) \quad \text{1}$$

is it correct to consider $\alpha(y, x) = 1$? Yes.

$$\text{If } \frac{\pi(y) q(y, x)}{\pi(x) q(x, y)} < 1 \Rightarrow \frac{\pi(x) q(x, y)}{\pi(y) q(y, x)} > 1 \Rightarrow \alpha(y, x) = \min(\dots, 1) = 1$$

$$2. \frac{\pi(y) q(y, x)}{\pi(x) q(x, y)} > 1$$

$$p(y, x) \pi(y) = q(y, x) \alpha(y, x) \pi(y) = q(y, x) \frac{\pi(x) q(x, y)}{\pi(y) q(y, x)} \stackrel{\text{if } \alpha > 1}{=} \pi(x) \cdot q(x, y) \quad \text{1}$$

This implies that $\pi(\cdot)$ is invariant:

$$\int_E P(x, A) \pi(x) dx = \pi(A) \quad \text{1}$$

12. Gaussian linear model with homoscedastic (iid) errors; conjugate prior when the variance is known/unknown, posterior and parameters updates; Jeffrey's prior (proportional to $1/\sigma^2$, where σ^2 is the error variance) and its posterior.

Suppose that we have the data (y_i, x_i) $i = 1, \dots, n$, where y_i are the responses and x_i are the covariates (p -dimensional). Suppose y_i to be continuous. We're interested in finding correlation between $\mathbb{E}[Y_i]$ and x_i . We assume:

$$Y | X, \beta, \sigma^2 \sim N_n(X\beta, \sigma^2 I_n) \quad X = n \times p \text{ matrix of covariates}$$

due to the modelization: $Y_i = X_i^T \beta + \varepsilon_i$ ($\varepsilon_i \sim N(0, \sigma^2)$).

The parameters of interest are (β, σ^2) or equivalently (β, τ) . (We consider σ^2 the variance and $\tau = 1/\sigma^2$ the precision).

The likelihood of the model is:

$$L(\beta, \sigma^2; y) \propto \frac{1}{(\sigma^2)^{n/2}} e^{-\frac{1}{2\sigma^2} (y - X\beta)^T (y - X\beta)}$$

Considering $\hat{\beta}_{MLE} = (X^T X)^{-1} X^T y$, $S^2 = (y - \hat{X}\beta_{MLE})^T (y - \hat{X}\beta_{MLE})$

$$\begin{aligned}
 (y - X\beta)^T (y - X\beta) &= (y - X\hat{\beta}_{MLE})^T (y - X\hat{\beta}_{MLE}) + (\beta - \hat{\beta}_{MLE})^T X^T X (\beta - \hat{\beta}_{MLE}) \\
 &\stackrel{\text{1}}{=} S^2 + (\beta - \hat{\beta}_{MLE})^T X^T X (\beta - \hat{\beta}_{MLE})
 \end{aligned}$$

$$\Rightarrow L(\beta, \sigma^2; y) \propto \frac{1}{(\sigma^2)^{n/2}} e^{-\frac{1}{2\sigma^2} [S^2 - (\beta - \hat{\beta}_{MLE})^T X^T X (\beta - \hat{\beta}_{MLE})]}$$

(12.) Conjugate prior when σ^2 is known.

Prior: $\beta \sim N_p(\underline{b}_0, B_0)$, B_0 invertible $p \times p$ matrix

$$\text{Posterior: } \beta|y \propto \exp\left\{-\frac{1}{2}(\beta - \hat{\beta}_{MLE})^T X^T X (\beta - \hat{\beta}_{MLE}) - \frac{1}{2}(\beta - \underline{b}_0)^T B_0^{-1}(\beta - \underline{b}_0)\right\}$$

$$\beta|y \sim N_p((X^T X + B_0^{-1})^{-1}(X^T X \hat{\beta}_{MLE} + B_0^{-1} \underline{b}_0), (X^T X + B_0^{-1})^{-1})$$

The posterior mean is a weighted average of the prior mean \underline{b}_0 and $\hat{\beta}_{MLE}$.

Conjugate prior when σ^2 is unknown.

We consider: $\pi(\beta, \sigma^2) = \pi(\beta|\sigma^2) \pi(\sigma^2)$

$$\text{Prior: } \begin{cases} \beta|\sigma^2 \sim N_p(\underline{b}_0, \sigma^2 B_0) \\ \sigma^2 \sim \text{inv-gamma}(\frac{v_0}{2}, \frac{J_0 \sigma_0^2}{2}) \end{cases} \quad \begin{matrix} \underline{b}_0 \in \mathbb{R}^p, B_0 \text{ invertible } p \times p \\ J_0, \sigma_0^2 > 0 \end{matrix}$$

Again we consider: $\pi(\beta, \sigma^2|y) = \pi(\beta|\sigma^2, y) \pi(\sigma^2|y)$

$$\text{Posterior: } \begin{cases} \beta|\sigma^2, y, X \sim N_p(\underline{b}_n, \sigma^2 B_n) \\ \sigma^2|y, X \sim \text{inv-gamma}(\frac{v_n}{2}, \frac{J_n \sigma_n^2}{2}) \end{cases}$$

$$B_n = (X^T X + B_0^{-1})^{-1}$$

$$\underline{b}_n = (X^T X + B_0^{-1})^{-1}(X^T X \hat{\beta}_{MLE} + B_0^{-1} \underline{b}_0)$$

$$J_n = J_0 + n$$

$$\sigma_n^2 = \frac{1}{J_n} (J_0 \sigma_0^2 + \underline{b}_0^T B_0^{-1} \underline{b}_0 + \underline{y}^T \underline{y} - \underline{b}_n^T B_n^{-1} \underline{b}_n)$$

Jeffrey's: prior: $\pi(\beta, \sigma^2) \propto \frac{1}{\sigma^2} \mathcal{U}(0, \infty) (\sigma^2)$

$$\text{posterior: } \begin{cases} \beta|\sigma^2, y, X \sim N_p(\hat{\beta}_{MLE}, \sigma^2(X^T X)^{-1}) \\ \sigma^2|y, X \sim \text{inv-gamma}(\frac{n-p}{2}, \frac{s^2}{2}) \end{cases} \quad \left. \begin{matrix} \text{proper if } X^T X \text{ is invertible} \\ \text{and } n > p \end{matrix} \right]$$

13. Formal approach to the model choice: calculation of the posterior probability of the model j , $j=1, \dots, K$.

In the formal approach of model selection we fix prior probability masses for any model. Then we use the Bayes theorem to compute the posterior distribution. The best models are those with highest posterior probability.

Let m be the index describing K models M_1, \dots, M_K .

- $P(m=j) =$ prior prob. to choose model M_j (typically $= 1/K$)
- model M_j : likelihood $f(y|\theta_j, M_j)$, prior $\pi(\theta_j|M_j) = \pi(\theta_j|m=j)$ where θ_j is the vector of parameters

How does it work?

1: for any model M_j compute the posterior for the parameters:
 $\pi(\theta_j|y, M_j) = \frac{f(y|\theta_j, M_j) \pi(\theta_j|M_j)}{m(y|M_j)}$, $m(y|M_j) = \int f(y|\theta_j, M_j) \pi(\theta_j|M_j) d\theta_j$,
 marginal distr. of data under model M_j

2: compute the posterior probability masses of M_1, \dots, M_K :

$$P(m=j|y) = \frac{m(y|M_j) P(m=j)}{m(y)}, \quad m(y) = \sum_{j=1}^K m(y|M_j) P(m=j)$$

3: choose the model with the highest $P(m=j|y)$. We can select 2/3 models and then compare these models through some goodness-of-fit criteria.

likelihood · prior
marginal

14. Computation of LPML (log-pseudo marginal likelihood) from a MCMC sample from the posterior, given all the datapoints.

Cross-validation:

We're going to split the dataset in training set and test set. The training we'll be used to compute the estimates, the test we'll be used for a comparison with the prediction computed through the training set.

Suppose to have n data. We split in $n-1$ and 1.

$y = (y_1, \dots, y_n) \Rightarrow y_{-i} = (y_1, \dots, y_{i-1}, y_{i+1}, \dots, y_n)$ training, y_i test
For all i ($i=1, \dots, n$) we find $\mathbb{E}(Y_i | Y_{-i})$ and compute

$$CPO_i := f_i(y_i | y_{-i}) = \text{CONDITIONAL PREDICTIVE ORDINATE } i$$

We suppose Y_1, \dots, Y_n s.t. $Y_i | \theta \sim f_i$;

$$\begin{aligned} f_i(y_i | y_{-i}) &= \int f_i(y_i | \theta) \pi(\theta | y_{-i}) d\theta \\ &= \int f_i(y_i | \theta) \frac{\pi(\theta) \prod_{j \neq i} f_j(y_j | \theta)}{\int \pi(\theta) \prod_{j \neq i} f_j(y_j | \theta) d\theta} d\theta \end{aligned}$$

= likelihood integrated w.r.t. the posterior distribution using all the data but the i -th

$$\begin{aligned} \frac{1}{f_i(y_i | y_{-i})} &= \frac{\int \pi(\theta) \prod_{j \neq i} f_j(y_j | \theta) d\theta}{\int \pi(\theta) \prod_{j=1}^n f_j(y_j | \theta) d\theta} = \frac{\int \frac{1}{\prod_{j=1}^n f_j(y_j | \theta)} \prod_{j \neq i} f_j(y_j | \theta) \pi(\theta) d\theta}{\int \pi(\theta) \prod_{j=1}^n f_j(y_j | \theta) d\theta} \\ &= \int \frac{1}{f_i(y_i | \theta)} \pi(\theta | y_1, \dots, y_n) d\theta \quad \text{posterior distribution using all the data} \\ &\approx \frac{1}{M} \sum_{m=1}^M \frac{1}{f_i(y_i | \theta^{(m)})} \quad (\theta^{(m)} = \text{MCMC sample}) \end{aligned}$$

$$\Rightarrow \hat{CPO}_i = \left(\frac{1}{M} \sum_{m=1}^M \frac{1}{f_i(y_i | \theta^{(m)})} \right)^{-1}$$

The higher the CPO_i , the better the model explains the i -th observation.

For model comparison we calculate the LOG-PSEUDO-MARGINAL-LIKELIHOOD

$$LPML := \sum_{i=1}^n \log(CPO_i) \quad (\text{the higher the better})$$

15. Probit regression model with latent variables: comp. of the full-conditionals

We consider the data $(y_i, x_i) \quad i=1, \dots, n, \quad y_i \in \{0, 1\}$.

We assume:

$$Y_i | X_i, \beta \sim \text{Be}(\pi_i), \quad \pi_i = \phi(x_i^\top \beta)$$

$$\beta \sim \pi(\beta) \quad (\text{commonly: } \beta \sim N_p(b_0, B_0))$$

This means that: $P(Y_i = 1) = \pi_i = \phi(x_i^\top \beta)$.

We introduce latent variables z_1, \dots, z_n :

$$z_i = x_i^\top \beta + \varepsilon_i, \quad \varepsilon_i \sim N(0, 1) \Rightarrow z_i \stackrel{iid}{\sim} N(x_i^\top \beta, 1)$$

and we consider:

$$Y_i = \begin{cases} 1 & \text{if } z_i > 0 \\ 0 & \text{otherwise} \end{cases}$$

$$1 - \phi(-k) =$$

The models are equivalent:

$$P(Y_i = 1) = P(z_i > 0) = 1 - P(z_i \leq 0) = 1 - P\left(\frac{z_i - x_i^\top \beta}{\sqrt{1}} \leq \frac{-x_i^\top \beta}{\sqrt{1}}\right) = \phi(x_i^\top \beta)$$

(25.) The parameters now are: (β, z) .

A Gibbs sampler can simply sample from $\pi(\beta, z|y)$, then we'll focus on the marginal of β . ($\pi(\beta|y)$)

$$\begin{aligned}\pi(\beta, z|y) &\propto \chi(\beta, z, Y) = \chi(Y|z) \chi(z|\beta) \chi(\beta) \\ &\propto \prod_{i=1}^n \underbrace{\left[(\mathbb{1}_{\{Y_i=1\}} \mathbb{1}_{\{z_i>0\}} + \mathbb{1}_{\{Y_i=0\}} \mathbb{1}_{\{z_i\leq 0\}}) \cdot N(z_i; x_i^T \beta, 1) \right]}_{\chi(Y_i|z)} \underbrace{\pi(\beta)}_{\chi(z|\beta)}\end{aligned}$$

The full-conditionals are:

- $\beta|z, y \propto \chi(Y, z, \beta) = \chi(Y|z, \beta) \chi(z|\beta) \chi(\beta)$

$\propto \chi(z|\beta) \pi(\beta)$. This last distribution is the posterior of a linear model with Gaussian likelihood where the data are z :

If $\pi(\beta) = \text{constant} \Rightarrow [\beta|z, y] \sim N_p((X^T X + B_0^{-1})^{-1} X^T y + B_0^{-1} b_0)$

If $\pi(\beta) = N_p(b_0, B_0) \rightarrow [\beta|z, y] \sim N_p(\tilde{\beta}, \tilde{B})$

where: $\begin{cases} \tilde{B} = (X^T X + B_0^{-1})^{-1} \\ \tilde{\beta} = \tilde{B}(X^T X \beta + B_0^{-1} b_0) \\ \hat{\beta} = (X^T X)^{-1} X^T y \end{cases}$

- $z| \beta, y \propto \chi(Y, z, \beta) = \chi(Y|z, \beta) \chi(z|\beta) \chi(\beta)$

$$\propto \prod_{i=1}^n \left[(\mathbb{1}_{\{Y_i=1\}} \mathbb{1}_{\{z_i>0\}} + \mathbb{1}_{\{Y_i=0\}} \mathbb{1}_{\{z_i\leq 0\}}) \cdot N(z_i; x_i^T \beta, 1) \right]$$

because of the product:

$$z_i | \beta, y_i = \begin{cases} \mathbb{1}_{\{z_i>0\}} \cdot N(z_i; x_i^T \beta, 1) & y_i = 1 \\ \mathbb{1}_{\{z_i\leq 0\}} \cdot N(z_i; x_i^T \beta, 1) & y_i = 0 \end{cases}$$

$$\Rightarrow \begin{cases} z_i | \beta, y_i = 1 \text{ is a } N(x_i^T \beta, 1) \text{ truncated to } z_i \in (0, \infty) \\ z_i | \beta, y_i = 0 \text{ is a } N(x_i^T \beta, 1) \text{ truncated to } z_i \in (-\infty, 0] \end{cases}$$

"truncated to $(0, \infty)$ ":



16. Likelihood for right-censored (conditionally independent) data

$$\begin{aligned} T_i &= \text{lifetime / failure time} & i = 1, \dots, n \\ C_i &= \text{time of censoring} & i = 1, \dots, n \end{aligned} \quad \left. \right\} \perp\!\!\!\perp$$

We have the data (y_i, δ_i) $i=1, \dots, n$: $y_i = \min(T_i, C_i)$, $\delta_i = \mathbb{1}_{\{T_i \leq C_i\}}$

We assume:

$T_i | \theta \sim f_i(t|\theta)$, $S_i = \text{survival function}$, $h_i = \text{hazard function}$

Then, the likelihood given the data is:

$$L(\theta | D) \propto \prod_{i=1}^n \left[(f_i(y_i|\theta))^{\delta_i} (S_i(y_i|\theta))^{1-\delta_i} \right]$$

which means:

- the contribution to the likelihood is:
- $f_i(y_i|\theta)$ if $\delta_i = 1$
 - $S_i(y_i|\theta)$ if $\delta_i = 0$

proof.

(Instead of $\text{IP}(T_i \in dy | \theta)$ we use $\text{IP}(T_i = y | \theta)$)

The contribution of the i -th observation is:

$$\begin{aligned} L(\theta | y_i, \delta_i) &= \left\{ \begin{array}{ll} \text{IP}(T_i = y_i, \delta_i = 1 | \theta) & \delta_i = 1 \\ \text{IP}(C_i = y_i, \delta_i = 0 | \theta) & \delta_i = 0 \end{array} \right. \\ &\stackrel{T_i \perp\!\!\!\perp C_i}{=} \left\{ \begin{array}{ll} \text{IP}(T_i = y_i, T_i \leq C_i | \theta) & \delta_i = 1 \\ \text{IP}(C_i = y_i, T_i > C_i | \theta) & \delta_i = 0 \end{array} \right. \\ &\stackrel{C_i \perp\!\!\!\perp \theta}{=} \left\{ \begin{array}{ll} \text{IP}(T_i = y_i, C_i \geq y_i | \theta) & \delta_i = 1 \\ \text{IP}(C_i = y_i, T_i > y_i | \theta) & \delta_i = 0 \end{array} \right. \\ &= \left\{ \begin{array}{ll} \text{IP}(T_i = y_i | \theta) \text{IP}(C_i \geq y_i | \theta) & \delta_i = 1 \\ \text{IP}(C_i = y_i | \theta) \text{IP}(T_i > y_i | \theta) & \delta_i = 0 \end{array} \right. \\ &= \left\{ \begin{array}{ll} f_i(y_i | \theta) \text{count}_1 & \delta_i = 1 \\ S_i(y_i | \theta) \text{count}_2 & \delta_i = 0 \end{array} \right. \end{aligned}$$

The following assumptions must hold:

1. T_i and C_i are independent (independent censoring)
2. The distribution of C_i does not depend on parameters contained in the distribution of T_i (non-informative censoring).
In our case this means $C_i \perp\!\!\!\perp \theta$
3. All the subjects in the sample are independent conditionally to all the parameters.

we proved what is the contribution to the likelihood in case of $\delta_i = 0$ or $\delta_i = 1$.
The (total) likelihood is the joint distribution of all the data.

17. First and second moments of the Dirichlet Process P (included the covariance).

Finite-dimensional Dirichlet distribution:

$$(D_1, \dots, D_{k-1}, D_k) \sim \text{Dir}(d_1, \dots, d_{k-1}, d_k) \quad d_j > 0, \quad D_j \in [0, 1] \quad \forall j, \quad \sum_j D_j = 1$$

If the distribution of (D_1, \dots, D_k) cannot be absolutely continuous w.r.t. the Lebesgue measure on \mathbb{R}^k , but the density of (D_1, \dots, D_{k-1}) is:

$$f(x_1, \dots, x_{k-1}) = \frac{\Gamma(d_1 + \dots + d_{k-1})}{\prod_{j=1}^{k-1} \Gamma(d_j)} \left[\prod_{j=1}^{k-1} x_j^{d_j-1} \right] (1-x_1 - \dots - x_{k-1})^{d_{k-1}-1} \mathbb{1}_{S_{k-1}}(x)$$

$$\text{where } S_{k-1} = \{(x_1, \dots, x_{k-1}) : \in \mathbb{R}^{k-1} : 0 \leq x_1 + \dots + x_{k-1} \leq 1, \quad 0 \leq x_j \leq 1 \quad \forall j = 1, \dots, k-1\}$$

Dirichlet Process:

Let α be a finite measure on \mathbb{R} : $0 < \alpha(\mathbb{R}) < \infty$ ($a := \alpha(\mathbb{R})$, $\alpha_0(A) = \frac{\alpha(A)}{a}$)
we say that P is a Dirichlet Process if for any finite and measurable partition of \mathbb{R} , say A_1, \dots, A_K , we have:

$$(P(A_1), \dots, P(A_K)) \sim \text{Dir}(\alpha(A_1), \dots, \alpha(A_K)) \quad : \quad P \sim DP(\alpha, \alpha_0)$$

Properties:

$$1. \quad P(A) \sim \text{Beta}(\alpha(A), \alpha(A^c)) \quad \forall A \in \mathcal{B}(\mathbb{R})$$

$$2. \quad E[P(A)] = \alpha_0(A), \quad \text{Var}(P(A)) = \frac{\alpha_0(A)\alpha_0(A^c)}{a+1} \quad \forall A \in \mathcal{B}(\mathbb{R})$$

Proof:

$$P(A) \sim \text{Beta}(\cdot, \cdot) \Rightarrow E[P(A)] = \frac{\alpha(A)}{\alpha(A) + \alpha(A^c)} = \frac{\alpha(A)}{\alpha(\mathbb{R})} = \alpha_0(A)$$

$$\text{Var}(P(A)) = \frac{\alpha(A)\alpha(A^c)}{(\alpha(A) + \alpha(A^c))^2 (\alpha(A) + \alpha(A^c) + 1)} = \frac{\alpha(A)\alpha(A^c)}{a^2(a+1)} = \frac{\alpha_0(A)\alpha_0(A^c)}{a+1}$$

$$3. \quad E[P(A)P(B)] = \frac{\alpha(A \cap B) + \alpha(A)\alpha(B)}{a(a+1)}$$

$$\text{Cov}(P(A), P(B)) = \frac{\alpha_0(A \cap B) - \alpha_0(A)\alpha_0(B)}{a+1}$$

$\forall A, B \in \mathcal{B}(\mathbb{R})$

Proof:

1. Assume $A \cap B = \emptyset$. We generate the partition $\{A, B, (A \cup B)^c\}$ of \mathbb{R} .

$$(P(A), P(B), P((A \cup B)^c)) \sim \text{Dir}(\alpha(A), \alpha(B), \alpha((A \cup B)^c))$$

and do we know the density of $(P(A), P(B))$: (finite-dim. Dir)

$$\begin{aligned} E[P(A)P(B)] &= \int_{S_2} x_1 x_2 \frac{\Gamma(a)}{\Gamma(\cdot)\Gamma(\cdot)\Gamma(\cdot)} x_1^{\alpha(A)-1} x_2^{\alpha(B)-1} (1-x_1-x_2)^{\alpha((A \cup B)^c)-1} dx_1 dx_2 \\ &= \frac{\alpha(A)}{\Gamma(a-1)\Gamma(1)\Gamma(a)} \cdot \frac{\Gamma(\alpha(A)+1)\Gamma(\alpha(B)+1)\Gamma(\alpha((A \cup B)^c))}{\Gamma(a+2)} \\ &\quad \frac{\alpha(A)\alpha(B)}{(a+1)a} \end{aligned}$$

$$\text{Cov}(P(A), P(B)) = E[P(A)P(B)] - E[P(A)]E[P(B)]$$

$$= \frac{\alpha(A)\alpha(B)}{(a+1)a} - \alpha_0(A)\alpha_0(B) = -\frac{\alpha_0(A)\alpha_0(B)}{a+1}$$

2. If $A \cap B \neq \emptyset$:

$$A = (A \cap B) \cup (A \setminus B)$$

$$\text{Cov}(P(A), P(B)) = \text{Cov}(P(A \cap B) + P(A \setminus B), P(A \cap B) + P(B \setminus A))$$

$$= \text{Var}(P(A \cap B)) + \underbrace{\text{Cov}(P(A \setminus B), P(A \cap B))}_{\text{disjoint, we are 1.}} + \dots$$

!! 18. Joint marginal distr. of a sample from a Dirichlet Process; its interpretation as the generalized Pólya urn.

$$X_1, \dots, X_n | P \sim p, \quad P \sim D_\alpha, \quad P | \alpha, \alpha_0 \sim DP(\alpha, \alpha_0)$$

We want to describe $\chi(X_1, \dots, X_n)$: $\chi(X_1, \dots, X_n) = \chi(X_1) \chi(X_2 | X_1) \dots \chi(X_n | X_1, \dots, X_{n-1})$

$$\begin{aligned} P(X_1 \in A) &= \int P(X_1 \in A | P) D_\alpha(dP) = \int \chi(X_1 | dP) = \int \chi(X_1 | P) \chi(dP) \\ &= \int_p P(A) D_\alpha(dP) = \mathbb{E}[P(A)] = \alpha_0(A) \end{aligned}$$

$$\rightarrow \chi(X_1) = \alpha_0$$

$$\begin{aligned} \chi(X_2 | X_1) &= \int_p \chi(X_2 | dP | X_1) = \int_p \chi(X_2 | P, X_1) \chi(dP | X_1) \\ &= \int_p P \cdot D_\alpha + \delta_{X_1}(dP) = \mathbb{E}[P | X_1] = \frac{\alpha + \delta_{X_1}}{\alpha+1} = \frac{\alpha}{\alpha+1} \alpha_0 + \frac{1}{\alpha+1} \delta_{X_1} \end{aligned}$$

Written differently: $X_2 | X_1 = \begin{cases} \sim \alpha_0 & \text{with prob. } \alpha/\alpha+1 \\ = X_1 & \text{with prob. } 1/\alpha+1 \end{cases}$

$$\begin{aligned} \chi(X_i | X_{i-1}, \dots, X_1) &= \int_p \chi(X_i | dP | X_{i-1}, \dots, X_1) = \int_p \chi(X_i | P, X_{i-1}, \dots, X_1) \chi(dP | X_{i-1}, \dots, X_1) \\ &= \int_p P \cdot D_\alpha + \sum_{j=1}^{i-1} \delta_{X_j}(dP) = \mathbb{E}[P | X_1, \dots, X_{i-1}] \\ &= \frac{\alpha + \sum_{j=1}^{i-1} \delta_{X_j}}{\alpha+i-1} = \frac{\alpha}{\alpha+i-1} \alpha_0 + \frac{1}{\alpha+i-1} \sum_{j=1}^{i-1} \delta_{X_j} \end{aligned}$$

Written differently:

$$X_i | X_{i-1}, \dots, X_1 = \begin{cases} \sim \alpha_0 & \text{with prob. } \alpha/\alpha+i-1 \\ = X_1 & \text{with prob. } 1/\alpha+i-1 \\ = X_2 & \text{with prob. } 1/\alpha+i-1 \\ \vdots & \vdots \\ = X_{i-1} & \text{with prob. } 1/\alpha+i-1 \end{cases}$$

$$\Rightarrow \chi(X_1, \dots, X_n) = \alpha_0 \prod_{i=2}^n \left(\frac{\alpha_0 + \sum_{j=1}^{i-1} \delta_{X_j}}{\alpha+i-1} \right) := \begin{array}{c} \text{GENERALIZED} \\ \text{PÓLYA URN} \\ (\text{Chinese Restaurant Process}) \end{array}$$

The generalized Pólya urn is a sampling scheme. We consider a urn. At each step we can sample a ball of the j -th color with probability proportional to the current number of balls of color j . We can also sample a ball of a new color from an other urn with an infinite number of balls $\alpha(\cdot)$, where each color is unique.

By sampling from a generalized Pólya urn we obtain a trajectory from a Dirichlet process.

Predictive for the Pólya urn: /for the Dirichlet Process

$$\chi(X_{n+1} | X_1, \dots, X_n) = \frac{\alpha}{\alpha+n} \alpha_0 + \frac{1}{\alpha+n} \sum_{j=1}^n \delta_{X_j}$$

The resulting distribution over labels with this sampling scheme is the same as the distribution over values in a Dirichlet process.

Stick-breaking scheme:

- α = finite measure on \mathbb{R}
- $\theta_1, \theta_2, \dots \stackrel{iid}{\sim} \alpha_0 (= \frac{\alpha(0)}{\alpha})$
- $Y_1, Y_2, \dots \stackrel{iid}{\sim} \text{Beta}(1, \alpha)$
- $V_1 = Y_1$
- $V_j = Y_j \prod_{i=1}^{j-1} (1 - Y_i) \quad j = 2, 3, \dots$

$$\left\{ \Rightarrow P := \sum_{j=1}^{+\infty} V_j \delta_{\theta_j} \sim D_\alpha \right.$$

Note: $\sum_{j=1}^{+\infty} V_j = 1$

19. Polya urn (with 2 colors): the sequence of random v.s. representing the color of the sampled balls is exchangeable with Beta distr. as Finetti measure.

Characterization of exchangeability:

$$X_1, \dots, X_n \text{ ir ex. iff } \exists Q \text{ prob. meas. : } P(X_1=x_1, \dots, X_n=x_n) = \prod_{i=1}^n P(X_i=x_i | \theta) Q(d\theta)$$

We consider a Polya urn:

we have a urn with balls of 2 colors, black and white. At each step we sample a ball with probability proportional to the number of balls of the same color currently in the urn. Once we sampled a ball, we return the ball in the urn and we add a new ball of the same color.

let $X_i = i\text{-th Bernoulli sampling}$: $X_i = \begin{cases} 1 & \text{if black} \\ 0 & \text{if white} \end{cases}$,

$W = \# \text{initial white balls}$, $B = \# \text{initial black balls}$, $S_n = \# \text{black balls at } n \text{ trial}$

A sequence is exchangeable:

$$P((1, 0)) = \frac{B}{B+W} \cdot \frac{W}{B+W+1} = P(0, 1)$$

$$\text{more generally: } P(X_1, \dots, X_n) = \frac{\Gamma(B+S_n) \Gamma(W+(n-S_n)) \Gamma(B+W)}{\Gamma(B) \Gamma(W) \Gamma(B+W+n)}$$

From de Finetti theorem: $P(X_1, \dots, X_n) = \prod_{i=1}^n \int \theta^{x_i} (1-\theta)^{1-x_i} Q(d\theta) = \int \theta^{S_n} (1-\theta)^{W-S_n} Q(d\theta)$

We call Q the de Finetti distr. $\lim_{n \rightarrow \infty} \bar{X}_n \sim \text{Beta}(B, W)$.

We have:

$$\begin{aligned} P(X_1, \dots, X_n) &= \int \theta^{S_n} (1-\theta)^{W-S_n} Q(d\theta) = \int \theta^{S_n} (1-\theta)^{W-S_n} \frac{\theta^{B-1} (1-\theta)^{W-1}}{B(B, W)} d\theta \\ &\stackrel{?}{=} \frac{B(B+S_n, W+(n-S_n))}{B(B, W)} \int \frac{\theta^{B+S_n-1} (1-\theta)^{W+n-S_n-1}}{B(B+S_n, W+(n-S_n))} d\theta \\ &\stackrel{?}{=} \frac{\Gamma(B+S_n) \Gamma(W+(n-S_n)) \Gamma(B+W)}{\Gamma(B) \Gamma(W) \Gamma(B+W+n)} = \text{probability associated to an exchangeable sequence produced from a Polya urn} \end{aligned}$$

20. Predictive distribution under the DPM model, marginalizing (i.e. integrating out) the r.p.m. P (Dirichlet Process)

Dirichlet Process Mixtures

We assume:

$$\left\{ \begin{array}{l} X_1, \dots, X_n | P \stackrel{\text{iid}}{\sim} f(x)(w) = \int_{\Theta} k(x; \theta) P(d\theta) \\ P \sim \text{D}_{\alpha}, \quad \alpha = \alpha \cdot \delta_0 \end{array} \right.$$

$$\Leftrightarrow \left\{ \begin{array}{l} X_i | \theta_i \stackrel{\text{iid}}{\sim} k(\cdot; \theta_i) \quad i=1, \dots, n \\ \theta_1, \dots, \theta_n | P \stackrel{\text{iid}}{\sim} p \\ p \sim \text{D}_{\alpha} \\ X_1, \dots, X_n | \theta \perp\!\!\!\perp p | \theta \end{array} \right.$$

where $k(x; \theta)$ is a parametric kernel: a density or function of x and a measurable function as function of θ .

Predictive distribution:

$$\begin{aligned} \mathbb{X}(X_{n+1} | X_1, \dots, X_n) &= \int_P \mathbb{X}(X_{n+1}, P | X_1, \dots, X_n) dP = \int_P \mathbb{X}(X_{n+1} | P, X_1, \dots, X_n) \mathbb{X}(P | X_1, \dots, X_n) dP \\ &= \int_P \mathbb{X}(X_{n+1} | P) \mathbb{X}(P | X_1, \dots, X_n) dP = \int_P \int_{\Theta} k(x; \theta) P(d\theta) \mathbb{X}(P | X_1, \dots, X_n) dP \\ &= \int_{\Theta} k(x; \theta) \int_P \mathbb{X}(P | X_1, \dots, X_n) P(d\theta) dP = \int_{\Theta} k(x; \theta) \mathbb{E}[P(d\theta) | X_1, \dots, X_n] \\ &= \int_{\Theta} k(x; \theta) \left[\frac{\alpha}{\alpha+n} \delta_0(d\theta) + \sum_{j=1}^n \frac{1}{\alpha+n} \delta_{x_j} \right] \\ &= \frac{\alpha}{\alpha+n} \int_{\Theta} k(x; \theta) \delta_0(d\theta) + \frac{1}{\alpha+n} \sum_{j=1}^n k(x; \theta_j) \end{aligned}$$

21. Computation of the full conditionals from the Polya-urn scheme in DPM.

$$\left\{ \begin{array}{l} X_1, \dots, X_n | P \stackrel{\text{iid}}{\sim} F(x; \theta) = \int k(x; \theta) P(d\theta) \\ P \sim \mathbb{D}_\alpha, \quad \alpha = a \cdot \lambda_0 \end{array} \right. \Leftrightarrow \left\{ \begin{array}{l} X_i | \theta_i \stackrel{\text{iid}}{\sim} k(\cdot; \theta_i) \\ \theta_1, \dots, \theta_n | P \stackrel{\text{iid}}{\sim} p \\ p \sim \mathbb{D}_\alpha \\ X_1, \dots, X_n | \theta \perp\!\!\!\perp P | \theta \end{array} \right.$$

$$P | X_1, \dots, X_n \sim \int \mathbb{D}_\alpha + \sum_{i=1}^n \delta_{\theta_i} H(d\theta_1, \dots, d\theta_n | X_1, \dots, X_n)$$

where:

$$H(d\theta_1, \dots, d\theta_n | X_1 = x_1, \dots, X_n = x_n) \propto \prod_{i=1}^n k(x_i; \theta_i) \alpha(d\theta_i) \prod_{i=2}^n \left[\alpha(d\theta_i) + \sum_{j=1}^{i-1} \delta_{\theta_j}(d\theta_i) \right]$$

22. Cluster estimates (definitions)

We're in a mixture model framework:

$$\left\{ \begin{array}{l} X_i | \theta_i \stackrel{\text{iid}}{\sim} k(\cdot; \theta_i) \\ \theta_1, \dots, \theta_n | P \stackrel{\text{iid}}{\sim} p \\ p \sim \mathbb{D}_\alpha, \quad \alpha = a \cdot \lambda_0 \end{array} \right.$$

$\theta_1, \dots, \theta_n$ may show ties \Rightarrow we obtain a latent partition g_n in $\theta_1, \dots, \theta_n$.

$g_n = \{A_1, \dots, A_k\}$ is a latent partition of $\{1, \dots, n\}$, where each A_j contains all the data indexes of the group j . The data are grouped as:

$$X_\ell \sim X_m \text{ (same group)} \Leftrightarrow \theta_\ell = \theta_m$$

Thanks to the ties we have a prior induced by the model. As clustering of the data we choose a summary of the posterior distribution of g_n , for instance the value of the random partition minimizing a posterior loss function $\ell(\cdot, \cdot)$:

$$\hat{g}_n^* = \underset{\substack{\text{optimal partition} \\ \text{when } g=\text{true partition}}}{\arg \min}_{\hat{g}_n \in \text{space of partitions}} E[\ell(g, \hat{g}_n) | X_1, \dots, X_n], \quad \ell(\cdot, \cdot) = \text{how much we're voting describing one part. in terms of the other (distance b/w. partitions)}$$

(23.) Binder's loss function.

Let's introduce S_1, \dots, S_n : $S_i = j$ if $X_i \in j\text{-th group}$.
(in this way $A_j = \{i : S_i = j\}$).

The Binder loss function penalizes with c_1 if two elements belong to the same group in the true partition but not in the considered one, with c_2 the opposite.

$$\ell_B(g, \hat{g}_n) = \sum_{i,j} [c_1 \mathbb{1}_{\{S_i=S_j\}} \mathbb{1}_{\{\hat{S}_i \neq \hat{S}_j\}} + c_2 \mathbb{1}_{\{S_i \neq S_j\}} \mathbb{1}_{\{\hat{S}_i = \hat{S}_j\}}].$$

It can be shown that the solution (opt.) is given by the partition which, in binary representation, minimizes the distance from the posterior similarity matrix M (where $M_{ij} = P(S_i = S_j | X_1, \dots, X_n)$) and \hat{S}^* :

$$\hat{S}^* = \arg \min_{\hat{S}} \left[\sum_{i,j} |\mathbb{1}_{\{S_i = S_j\}} - M_{ij}| \right]$$

Ex. $B = \{\{1, 2, 3\}, \{12, 3\}, \{1, 23\}, \{13, 2\}, \{123\}\}$

Binary rep. $A_1 = \begin{bmatrix} 1 & 0 & 0 \\ 0 & 1 & 0 \\ 0 & 0 & 1 \end{bmatrix}$, $A_2 = \begin{bmatrix} 1 & 1 & 0 \\ 0 & 0 & 1 \\ 0 & 0 & 0 \end{bmatrix}$, $A_3 = \begin{bmatrix} 1 & 0 & 1 \\ 0 & 1 & 0 \\ 0 & 0 & 1 \end{bmatrix}$, $A_4 = \begin{bmatrix} 0 & 1 & 1 \\ 1 & 0 & 1 \\ 0 & 1 & 0 \end{bmatrix}$, $A_5 = \begin{bmatrix} 1 & 1 & 1 \\ 1 & 1 & 1 \\ 1 & 1 & 1 \end{bmatrix}$

and $M = \begin{bmatrix} 1 & 0.8 & 1 & 0.2 & 1 \\ 0.8 & 1 & 0.2 & 1 & 0.2 \\ 1 & 0.2 & 1 & 0.2 & 1 \\ 0.2 & 1 & 0.2 & 1 & 0.2 \\ 1 & 0.2 & 1 & 0.2 & 1 \end{bmatrix}$ \rightarrow posterior prob. that 1 and 2 are grouped together

$$d(A_1, M) = 0.8 + 0.2 + 0.2 = 1.2$$

$$d(A_2, M) = 0.6$$

$$d(A_3, M) = 1.8$$

$$d(A_4, M) = 1.8$$

$$d(A_5, M) = 2.9$$

\rightarrow optimal = A_2