

$$\begin{aligned} \underline{x} &= (x_1, \dots, x_n)^T \\ \underline{y} &= (y_1, \dots, y_n)^T \end{aligned} \quad \left. \begin{array}{l} \text{continuous } n\text{-dimensional} \\ \text{random vectors} \end{array} \right\}$$

We suppose $\exists g(\cdot)$ invertible s.t. \exists the partial derivatives and the Jacobian matrix has determinant $\neq 0$.

$\underline{x} \in \mathbb{R}^n$ random vector

$f_{\underline{x}}(\underline{x})$ is the density function that describes the distribution of \underline{x}

We express the density function of \underline{y} :

$$h_{\underline{y}}(\underline{y}) = f_{\underline{x}}(g^{-1}(\underline{y})) \cdot \|Jg^{-1}\|$$

absolute value of the determinant of the Jacobian of the inverse of g

$$Jg^{-1} = \begin{bmatrix} \frac{\partial x_1}{\partial y_1} & \dots & \frac{\partial x_1}{\partial y_n} \\ \vdots & & \vdots \\ \frac{\partial x_n}{\partial y_1} & \dots & \frac{\partial x_n}{\partial y_n} \end{bmatrix}$$

We're defining a new random variable by applying some tranfor.
(new n -dim rand.-var.)

GAMMA RANDOM VARIABLES

Gamma function:

$$\Gamma(\alpha) = \int_0^\infty t^{\alpha-1} e^{-t} dt$$

properties:

- $\Gamma(1) = 1$
- $\Gamma(\frac{1}{2}) = \sqrt{\pi}$

• $\boxed{\Gamma(\alpha+1) = \alpha \Gamma(\alpha)}$ $\forall \alpha > 0$

• $\Gamma(n+1) = n!$ $\forall n \in \{0, 1, \dots\}$

proof. • $\Gamma(\alpha+1) = \int_0^\infty x^\alpha e^{-x} dx = \int_0^\infty x^\alpha \frac{d}{dx} (-e^{-x}) dx$

$$= \left[-x^\alpha e^{-x} \right]_0^\infty + \alpha \int_0^\infty x^{\alpha-1} e^{-x} dx$$

$$= \alpha \Gamma(\alpha)$$

We want to define a random variable which density is proportional to: $t^{\alpha-1} e^{-t}$.
 X random variable s.t. $f_X(x) \propto x^{\alpha-1} e^{-x} \mathbf{1}_{(0,\infty)}(x)$

$$\begin{aligned} f_X(x) &= \frac{x^{\alpha-1} e^{-x} \mathbf{1}_{(0,\infty)}(x)}{\int_0^\infty t^{\alpha-1} e^{-t} dt} \\ &= \frac{1}{\Gamma(\alpha)} x^{\alpha-1} e^{-x} \mathbf{1}_{(0,\infty)}(x) \end{aligned} \quad \Rightarrow \quad X \sim \text{Gamma}(\alpha, 1) \quad \alpha \in \mathbb{R}^+$$

random var. taking values in $(0, \infty)$

We want to introduce another parameter to governate the shape.
 We'll introduce it by a transformation of $X \sim \text{Gamma}(\alpha, 1)$.

$$U := \frac{X}{\beta}, \quad \beta > 0$$

$$X = \beta U \Rightarrow dx = \beta du$$

The density of U becomes:

$$\begin{aligned} f_U(u) &= f_X(\beta u) \cdot \beta \mathbb{1}_{(0, \infty)}(u) \\ &\stackrel{\text{def}}{=} \frac{1}{\Gamma(\alpha)} (\beta u)^{\alpha-1} e^{-\beta u} \cdot \beta \mathbb{1}_{(0, \infty)}(u) \\ &= \frac{1}{\Gamma(\alpha)} \beta^\alpha u^{\alpha-1} e^{-\beta u} \mathbb{1}_{(0, \infty)}(u) \end{aligned} \implies U \sim \text{Gamma}(\alpha, \beta)$$

$\alpha = \text{shape parameter}$

$\beta = \text{rate parameter}$

$1/\beta = \text{scale parameter}$

$[-] e^{-\frac{u}{\beta}}$

To study the impact of the parameters:

$$\begin{aligned} \mathbb{E}[U^n] &= \int_{\mathbb{R}} u^n \frac{\beta^\alpha}{\Gamma(\alpha)} u^{\alpha-1} e^{-\beta u} \mathbb{1}_{(0, \infty)}(u) du \\ &\stackrel{\text{def}}{=} \frac{\beta^\alpha}{\Gamma(\alpha)} \int_0^\infty u^{\alpha+n-1} e^{-\beta u} du \\ &\stackrel{\text{kernel of a Gamma random var: } \text{Gamma}(\alpha+n, \beta)}{=} \int_0^\infty \frac{\beta^{\alpha+n}}{\Gamma(\alpha+n)} u^{\alpha+n-1} e^{-\beta u} du \\ &= \frac{1}{\beta^n} \frac{\Gamma(\alpha+n)}{\Gamma(\alpha)} = \frac{1}{\beta^n} \cdot \boxed{(n+\alpha-1)(n+\alpha-2) \cdots \alpha} = \frac{(\alpha)_n}{\beta^n} \\ &\quad := \text{rising factorial in } n \end{aligned}$$

RISING FACTORIAL OF a IN b :

$$\frac{\Gamma(a+b)}{\Gamma(a)} = (a)_b$$

$=$ truncation of a factorial term

$$\implies \mathbb{E}[U] = \frac{\alpha}{\beta}$$

$$\text{Var}(U) = \mathbb{E}[U^2] - \mathbb{E}^2[U] = \frac{\alpha}{\beta^2}$$

Special cases of the Gamma distribution:

$$1. \text{Gamma}(1, \beta) \sim \mathcal{E}(\beta)$$

$$2. \text{Gamma}\left(\frac{n}{2}, \frac{1}{2}\right) \sim \chi^2(n)$$

$$3. U \sim \text{Gamma}(n, \beta) \implies 2\beta U \sim \text{Gamma}\left(\frac{2n}{2}, \frac{1}{2}\right) \sim \chi^2(2n) \quad !(\text{E})$$

$$4. Z = U^{-1}, U \sim \text{Gamma}(\alpha, \beta) \implies Z \sim \text{inverse-gamma}(\alpha, \beta)$$

Properties:

$$1. \text{Additivity: } U_1, \dots, U_k \stackrel{\text{iid}}{\sim} \text{Gamma}(\alpha_j, \beta) \quad j=1, \dots, k \quad (U_j \sim \text{Gamma}(\alpha_j, \beta_j))$$

$$\implies \sum_{j=1}^k U_j \sim \text{Gamma}\left(\sum_{j=1}^k \alpha_j, \beta\right)$$

$$2. \left. \begin{array}{l} U_1 \sim \text{Gamma}(\alpha_1, \beta) \\ U_2 \sim \text{Gamma}(\alpha_2, \beta) \end{array} \right\} \implies \left\{ \begin{array}{l} (U_1 + U_2) \perp\!\!\!\perp \frac{U_1}{U_1 + U_2} \\ (U_1 + U_2) \perp\!\!\!\perp \frac{U_2}{U_1 + U_2} \end{array} \right.$$

BETA RANDOM VARIABLES

Beta function:

$$\beta(\alpha_1, \alpha_2) = \int_0^1 t^{\alpha_1-1} (1-t)^{\alpha_2-1} dt = \frac{\Gamma(\alpha_1) \Gamma(\alpha_2)}{\Gamma(\alpha_1 + \alpha_2)}$$

Similarly to the previous case:

We want to define a random var. R :

$$f_R(x) \propto x^{\alpha_1-1} (1-x)^{\alpha_2-1} \mathbb{1}_{(0,1)}(x)$$

$$\begin{aligned} \Rightarrow f_R(x) &= \frac{x^{\alpha_1-1} (1-x)^{\alpha_2-1}}{\int_0^1 t^{\alpha_1-1} (1-t)^{\alpha_2-1} dt} \mathbb{1}_{(0,1)}(x) \\ &= \frac{1}{\beta(\alpha_1, \alpha_2)} x^{\alpha_1-1} (1-x)^{\alpha_2-1} \mathbb{1}_{(0,1)}(x) \quad \Rightarrow R \sim \text{Beta}(\alpha_1, \alpha_2) \end{aligned}$$

random var. taking values in $(0, 1)$

Property: • $\begin{cases} U_1 \sim \text{Gamma}(\alpha_1, \beta) \\ U_2 \sim \text{Gamma}(\alpha_2, \beta) \end{cases} \Rightarrow \begin{cases} \frac{U_1}{U_1 + U_2} \sim \text{Beta}(\alpha_1, \alpha_2) \\ \frac{U_2}{U_1 + U_2} \sim \text{Beta}(\alpha_2, \alpha_1) \end{cases}$

Again: n -th moment:

$$\begin{aligned} \mathbb{E}[R^n] &= \int r^n \frac{1}{\beta(\alpha_1, \alpha_2)} r^{\alpha_1-1} (1-r)^{\alpha_2-1} \mathbb{1}_{(0,1)}(r) dr \\ &= \frac{1}{\beta(\alpha_1, \alpha_2)} \int_0^1 r^{\alpha_1+n-1} (1-r)^{\alpha_2-1} dr \\ &= \frac{\beta(\alpha_1+n, \alpha_2)}{\beta(\alpha_1, \alpha_2)} = \frac{(\alpha_1)_n}{(\alpha_1+\alpha_2)_n} = \frac{(\alpha_1+n-1)(\alpha_1+n-2)\dots\alpha_1}{(\alpha_1+\alpha_2+n-1)(\alpha_1+\alpha_2+n-2)\dots(\alpha_1+\alpha_2)} \end{aligned}$$

$$\Rightarrow \mathbb{E}[R] = \frac{\alpha_1}{\alpha_1 + \alpha_2}$$

$$\text{var}(R) = \mathbb{E}[R^2] - \mathbb{E}^2[R] = \frac{\alpha_1 \alpha_2}{(\alpha_1 + \alpha_2 + 1)(\alpha_1 + \alpha_2)^2}$$

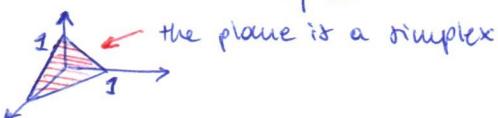
Special case:

$$1. \text{ Beta}(1,1) \sim \mathcal{U}([0,1])$$

Properties:

$$1. X \sim \text{Beta}(a,b) \implies 1-X \sim \text{Beta}(b,a)$$

$$2. \text{ If } \alpha_1 > 1, \alpha_2 > 1 \implies \text{mode} = \frac{\alpha_2-1}{\alpha_1+\alpha_2-1}$$

Betas are useful to model probabilities (since it takes values in $(0,1)$).If we want to model more than one outcome (more than one probability) the Dirichlet distribution is one of the possible multivariate extensions of the Beta. It takes values in the simplex:

DIRICHLET DISTRIBUTION

We start from considering more than 2 gammas.

$$U_1, \dots, U_k \stackrel{\text{iid}}{\sim} \text{Gamma}(\alpha_j, \beta) \quad (U_j \sim \text{Gamma}(\alpha_j, \beta))$$

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

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

We impose: $\sum_{j=1}^k U_j = 1$ (and so (X_1, \dots, X_k) takes values in:

$$\Rightarrow (X_1, \dots, X_k) \in \Delta_{k-1} = \left\{ x \in \mathbb{R}^k : \sum_{j=1}^k x_j = 1, x_j \in [0, 1] \right\}$$

$k-1$ simplex



Density of a Dirichlet distribution:

$$f_X(x) = \frac{1}{B(\alpha)} \prod_{j=1}^k x_j^{\alpha_j - 1} \mathbb{1}_{\Delta_{k-1}}(x) \quad B(\alpha) = \frac{\Gamma(\alpha_1) \cdots \Gamma(\alpha_k)}{\Gamma(\alpha_1 + \dots + \alpha_k)}$$

$$= \frac{1}{B(\alpha)} \left[\prod_{j=1}^{k-1} x_j^{\alpha_j - 1} \right] \left(1 - \sum_{j=1}^{k-1} x_j \right)^{\alpha_k - 1} \mathbb{1}_{\Delta_{k-1}}(x)$$

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

Property:

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

$$\text{because } X_j = \frac{U_j}{U_1 + \dots + U_k} = \frac{U_j}{U_j + \sum_{l=j+1}^k U_l}$$

$$\text{Var}(X_j) = \frac{\alpha_j(\alpha_0 - \alpha_j)}{\alpha_0^2 (\alpha_0 + 1)}$$

and so we're in the case of gammas (we have the property)

$$\text{Cov}(X_j, X_l) = -\frac{\alpha_l \alpha_j}{\alpha_0^2 (\alpha_0 + 1)} \quad j \neq l$$

(Other) properties:

1. The marginal distribution is: $X_j \sim \text{Beta}(\alpha_j, \alpha_0 - \alpha_j)$

2. Aggregation property:

$$(X_1, \dots, X_n) \sim \text{Dirichlet}(\alpha_1, \dots, \alpha_n)$$

$$(Z_1, \dots, Z_m) := \left(\sum_{i \in C_1} X_i, \dots, \sum_{i \in C_m} X_i \right)$$

$$\text{where } \begin{cases} C_i \subseteq \{1, \dots, n\} & \forall i \\ C_i \cap C_j = \emptyset & \forall i \neq j \\ \bigcup_{i=1}^m C_i = \{1, \dots, n\} \end{cases}$$

we sum up some dimensions
($m \leq n$)

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

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

$$3. (X_1, \dots, X_k) \sim \text{Dir}(\alpha_1, \dots, \alpha_k) \Rightarrow \left\{ \begin{array}{l} X_1 \sim \text{Beta}(\alpha_1, \alpha_0 - \alpha_1) \\ X_1 \perp \left(\frac{X_2}{X_1}, \dots, \frac{X_k}{X_1} \right) \end{array} \right.$$

$$\text{we can take out one dimension} \rightarrow \left[\left(\frac{X_2}{1-X_1}, \dots, \frac{X_k}{1-X_1} \right) \sim \text{Dir}(\alpha_2, \dots, \alpha_k) \right]$$

proof that we can take out a dimension:

We do a transformation:

18/09/2020

3/3

$$g(x_1, \dots, x_{k-1}) = \begin{cases} Q_1 = x_1 \\ Q_2 = \frac{x_2}{1-x_1} \\ \vdots \\ Q_{k-1} = \frac{x_{k-1}}{1-x_1} \end{cases} \begin{cases} \in (0,1) \\ \in \Delta_{k-2} \end{cases}$$

We have that $Q_k = \sum_{j=2}^{k-1} Q_j$.

$$g^{-1}(Q_1, \dots, Q_{k-1}) = \begin{cases} x_1 = Q_1 \\ x_2 = Q_2(1-Q_1) \\ \vdots \\ x_{k-1} = Q_{k-1}(1-Q_1) \end{cases}$$

$$J_{g^{-1}} = \begin{bmatrix} \frac{\partial x_1}{\partial Q_1} & \dots & \frac{\partial x_1}{\partial Q_{k-1}} \\ \vdots & & \vdots \\ \frac{\partial x_{k-1}}{\partial Q_1} & \dots & \frac{\partial x_{k-1}}{\partial Q_{k-1}} \end{bmatrix} = \begin{bmatrix} 1 & 0 & \dots & 0 \\ -Q_2 & 1-Q_1 & \dots & 0 \\ \vdots & & 1-Q_1 & \dots \\ 0 & 0 & \dots & 1-Q_1 \end{bmatrix}$$

$$\|J_{g^{-1}}\| = 1 \cdot (1-Q_1)^{k-2}$$

Starting from the density of X we want to derive the density of Q :

$$\begin{aligned} f_Q(q) &= f_X(g^{-1}(q)) \cdot \|J_{g^{-1}}\| \\ &= \frac{1}{\beta(\alpha)} \left[\prod_{j=2}^{k-1} (q_j(1-q_1))^{\alpha_j-1} \right] \cdot q_1^{\alpha_1-1} \cdot \left(1 - \sum_{j=2}^{k-1} q_j(1-q_1) - q_1 \right)^{\alpha_{k-1}-1} \cdot (1-q_1)^{k-2} \\ &= \frac{1}{\beta(\alpha)} q_1^{\alpha_1-1} (1-q_1)^{\sum_{j=2}^{k-1} \alpha_j} (1-q_1)^{\alpha_{k-1}-1} \left[\prod_{j=2}^{k-1} q_j \right] \left[1 - \sum_{j=2}^{k-1} q_j \right]^{\alpha_{k-1}-1} (1-q_1)^{k-2} \\ &= \frac{1}{\beta(\alpha)} q_1^{\alpha_1-1} (1-q_1)^{\sum_{j=1}^{k-1} \alpha_j-1} (1-q_1)^{k-2} \left[\prod_{j=2}^{k-1} q_j \right] \left[1 - \sum_{j=2}^{k-1} q_j \right]^{\alpha_{k-1}-1} \\ &= \frac{1}{\beta(\alpha)} q_1^{\alpha_1-1} (1-q_1)^{\alpha_0-\alpha_1-1} \left[\prod_{j=2}^{k-1} q_j \right] \left[1 - \sum_{j=2}^{k-1} q_j \right]^{\alpha_{k-1}-1} \end{aligned}$$

Beta random var.
up to a normalization
constant

Beta($\alpha_1, \alpha_0 - \alpha_1$)

Dirichlet distribution
defined on the $k-2$ -dimensional
simplex: the parameter α_1
is not contributing in any way
to the shape of the distribution

Let's now see some discrete distributions.

BERNOULLI DISTRIBUTION

$X \in \{0, 1\}$ describing the success/failure.

We want to model the randomness behind the trial.

We consider $\theta \in (0, 1)$; (θ random var.)

$$p_X(x) = \theta^x (1-\theta)^{1-x} \quad (\theta = \text{P(success)}) \implies X \sim \text{Be}(\theta)$$

probability mass function

$$\mathbb{E}[X] = \theta$$

$$\text{Var}(X) = \theta(1-\theta)$$

More than one trial:

BINOMIAL DISTRIBUTION

$X \in \{0, 1, \dots, n\}$, $\theta \in (0, 1)$

$$p_X(x) = \binom{n}{x} \theta^x (1-\theta)^{n-x} \implies X \sim \text{Bi}(n, \theta)$$

$$\mathbb{E}[X] = n\theta$$

$$\text{Var}(X) = n\theta(1-\theta)$$

(it describes the outcome of a sequence of trials)

MULTINOMIAL DISTRIBUTION

multivariate extension
of the Binomial distribution

\underline{R} = k-dimensional random vector characterized by: $\underline{\lambda} = (\lambda_1, \dots, \lambda_k)$

where λ_j = probability of success for the j-th dimension

$$\sum_{j=1}^k \lambda_j = 1, \quad \underline{\lambda} \in \Delta_{k-1}$$

$$p_{\underline{R}}(\underline{r}) = \frac{m!}{r_1! \cdots r_k!} \lambda_1^{r_1} \cdots \lambda_k^{r_k} \mathbb{1}_{\left\{ \begin{array}{l} \sum r_j = m, \\ r_j \in \{0, \dots, m\} \end{array} \right\}}$$

r_j = # success in the j-th dimension

m = # trials (total)

$$\mathbb{E}[R_j] = m \cdot \lambda_j$$

$$\text{Var}(R_j) = m \lambda_j (1-\lambda_j)$$

$$\text{Cov}(R_i, R_j) = -m \lambda_i \lambda_j \quad i \neq j$$

Bayes theorem:

$$A, B \text{ events}, \quad P(B) > 0 \quad : \quad P(A|B) = \frac{P(B|A) P(A)}{P(B)}$$

Intuitively: $A, B \in \Omega$:



In the case of probability mass functions:

$$\begin{aligned} p_{X|Y}(x|y) &= \frac{p_{Y|X}(y|x) p_X(x)}{p_Y(y)} = \frac{p_{Y|X}(y|x) p_X(x)}{\sum_x p_{X,Y}(x,y)} \\ &= \frac{p_{Y|X}(y|x) p_X(x)}{m(y)} \propto p_{X,Y}(x,y) \end{aligned}$$

we can update our beliefs just working proportionally to the joint distribution

Similarly, with density functions:

$$\begin{aligned} f_{X|Y}(x|y) &= \frac{f_{Y|X}(y|x) f_X(x)}{f_Y(y)} = \frac{f_{Y|X}(y|x) f_X(x)}{\int f_{X,Y}(x,y) dx} \\ &= \frac{f_{Y|X}(y|x) f_X(x)}{m(y)} \end{aligned}$$

The important thing is that we can use our prior belief to update the posterior (together with the data)

Remark: $P(B|A_1, A_2) \propto P(A_1, A_2|B) P(B)$

$$\propto P(A_2|B, A_1) P(A_1|B) P(B)$$

we can sequentially update our beliefs

Def. Conjugacy: The distribution of X is conjugate to the distribution of Y if $Y|X$ is in the same family of distributions of X .
(when we update our prior belief we stay in the same family of distributions)

BAYESIAN POINT ESTIMATION

We denote by $\Theta = \text{set of parameters} = \text{set of all the possible states } \theta, \theta \sim \pi(\theta)$.

We have a set of possible actions a : $a \in A$ is a single action.

We define a loss function: $l: \Theta \times A \rightarrow \mathbb{R}, l: (\theta, a) \mapsto l(\theta, a)$.

meaning: when the state is θ , how much do we lose choosing the action a ?

We work with:

$$E_\pi[l(\theta, a)] = \text{prior expected loss.}$$

The posterior expected loss will be:

$$E_\pi[l(\theta, a)|X]$$

where: $\left\{ \begin{array}{l} X_1, \dots, X_n | \theta \sim f(x, \theta) \\ \theta \sim \pi(\theta) \end{array} \right. \text{ s.t. } \theta \in \Theta \subseteq \mathbb{R}^k$

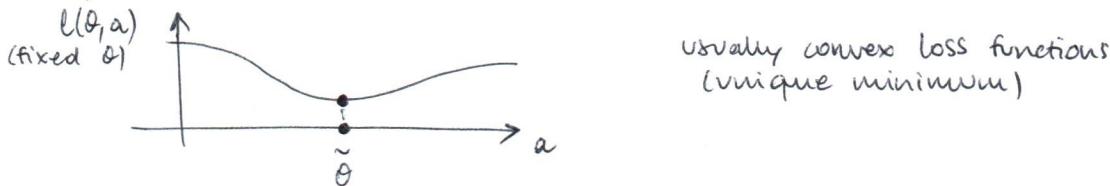
So, we choose an action a to which corresponds a loss. Due to the fact that θ (based on which we choose a) is random, instead of minimizing $l(\theta, a)$ we have to take care also of the distribution of θ . And so, we work with prior and posterior expectation of $l(\theta, a)$ w.r.t. π .

A way to choose an optimal θ (in $\pi(\theta)$) is to set $a = \hat{\theta}$:

$$\Rightarrow \underset{a \in \Theta}{\operatorname{arg\,min}} \mathbb{E}_{\pi}[l(\theta, a) | X] = \hat{\theta}$$

We can choose the loss function in several ways.

A good choice is a convex function:



Different loss functions lead to different point estimations.

Exercise 1. (Quadratic loss function leads to posterior expectation)

$$\theta \in \Theta \subseteq \mathbb{R}, \quad a \in \Theta, \quad l(\theta, a) := (\theta - a)^2$$

If we assume this loss function:

$$\Rightarrow \hat{\theta} = \mathbb{E}_{\pi}[\theta | X] = \underset{a \in \Theta}{\operatorname{arg\,min}} \mathbb{E}_{\pi}[l(\theta, a) | X]$$

↑
we have to prove it

$$\begin{aligned} \mathbb{E}_{\pi}[l(\theta, a) | X] &= \int_{\Theta} l(\theta, a) \pi(\theta | X) d\theta & \pi(\theta | X) &= \text{posterior distribution} \\ &= \int_{\Theta} (\theta - a)^2 \pi(\theta | X) d\theta \\ &= \int_{\Theta} (\theta^2 - 2a\theta + a^2) \pi(\theta | X) d\theta \end{aligned}$$

$$\frac{d}{da} (\theta^2 - 2a\theta + a^2) = 2a - 2\theta$$

We have that the minimum is given by:

$$\begin{aligned} \int_{\Theta} (2a - 2\theta) \pi(\theta | X) d\theta &= 0 \quad \Rightarrow \quad 2a \underbrace{\int_{\Theta} \pi(\theta | X) d\theta}_1 = 2 \underbrace{\int_{\Theta} \theta \pi(\theta | X) d\theta}_{\mathbb{E}[\theta | X]} \\ \Rightarrow a = \hat{\theta} = \mathbb{E}_{\pi}[\theta | X] &= \text{the point estimate} \\ &\quad \text{is the posterior mean} \end{aligned}$$

Exercise 2. (Quadratic loss function leads to posterior expectation (higher dimension))

Same as before but: $l(\underline{\theta}, \underline{a}) = (\underline{\theta} - \underline{a})^T Q (\underline{\theta} - \underline{a})$, $Q = \text{symmetric positive semidefinite matrix}$

(we have higher dimensions)

$$\underline{\theta} \in \Theta \subseteq \mathbb{R}^P, \quad \underline{a} \in \Theta$$

$$\Rightarrow \hat{\underline{\theta}} = \mathbb{E}_{\pi}[\underline{\theta} | X] = \underset{\underline{a} \in \Theta}{\operatorname{arg\,min}} \mathbb{E}_{\pi}[l(\underline{\theta}, \underline{a}) | X]$$

we want to prove this

We directly start from:

$$\begin{aligned} \frac{d}{d\underline{a}} \mathbb{E}_{\pi}[l(\underline{\theta}, \underline{a}) | X] &= \frac{d}{d\underline{a}} \int_{\Theta} (\underline{\theta} - \underline{a})^T Q (\underline{\theta} - \underline{a}) \pi(\underline{\theta} | X) d\underline{\theta} \\ &:= F(\underline{a}) \end{aligned}$$

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

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

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

$$\Rightarrow -2Q \underbrace{\int_{\Theta} \underline{\theta} \pi(\underline{\theta} | \underline{x}) d\underline{\theta}}_{E_{\pi}[\underline{\theta} | \underline{x}]} = -2Q \int_{\Theta} \underline{\theta} \pi(\underline{\theta} | \underline{x}) d\underline{\theta}$$

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

Exercise 3. (Linear loss function leads to quantiles)

$$\Theta \in \mathbb{R} \subseteq \mathbb{R}, \quad \mathcal{U} = \mathbb{R}, \quad a \uparrow \quad : \quad \ell(\theta, a) = \begin{cases} k_0(\theta - a) & a < \theta \\ k_1(a - \theta) & a \geq \theta \end{cases}$$

With this loss function our point estimate will be:

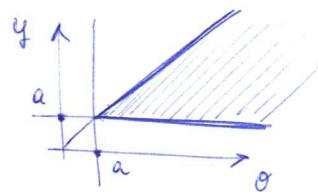
$$\tilde{\theta} = \text{quantile of } \pi(\theta | \underline{x}) \text{ of order } \frac{k_0}{k_1 + k_0} \quad \leftarrow \text{we want to prove this}$$

We start from:

$$\begin{aligned} E_{\pi}[\ell(\theta, a) | \underline{x}] &= \int_{\Theta} \ell(\theta, a) \pi(\theta | \underline{x}) d\theta \\ &= \int_{\Theta} k_0(\theta - a) \mathbf{1}_{(a < \theta)} \pi(\theta | \underline{x}) d\theta + \int_{\Theta} k_1(a - \theta) \mathbf{1}_{(a \geq \theta)} \pi(\theta | \underline{x}) d\theta \\ &= \int_a^{\infty} k_0(\theta - a) \pi(\theta | \underline{x}) d\theta + \int_{-\infty}^a k_1(a - \theta) \pi(\theta | \underline{x}) d\theta \\ &= k_0 I_0 + k_1 I_1 \end{aligned}$$

$$\begin{aligned} I_0 &= \int_a^{\infty} (a - \theta) \pi(\theta | \underline{x}) d\theta = \int_a^{\infty} \int_a^{\theta} dy \pi(\theta | \underline{x}) d\theta \\ &= \int_a^{\infty} \int_y^{\infty} \pi(\theta | \underline{x}) d\theta dy \\ &= \int_a^{\infty} \mathbb{P}(\theta > y | \underline{x}) dy \end{aligned}$$

$$\begin{aligned} I_1 &= \int_{-\infty}^a (a - \theta) \pi(\theta | \underline{x}) d\theta = \int_{-\infty}^a \int_{\theta}^a dy \pi(\theta | \underline{x}) d\theta \\ &= [\dots] = \int_{-\infty}^a \mathbb{P}(\theta < y | \underline{x}) dy \end{aligned}$$



$$\Rightarrow E_{\pi}[\ell(\theta, a) | \underline{x}] = k_0 \int_a^{\infty} P(\theta > y | \underline{x}) dy + k_1 \int_{-\infty}^a P(\theta < y | \underline{x}) dy$$

$$\begin{aligned} \frac{d}{da} E_{\pi}[\ell(\theta, a) | \underline{x}] &= -k_0 P(\theta > a | \underline{x}) + k_1 P(\theta < a | \underline{x}) \\ &\stackrel{!}{=} -k_0 (1 - P(\theta \leq a | \underline{x})) + k_1 P(\theta < a | \underline{x}) = 0 \end{aligned}$$

$$\Rightarrow P(\theta < a | \underline{x}) = \frac{k_0}{k_1 + k_0} \quad \Rightarrow \quad P(\theta < \tilde{\theta} | \underline{x}) = \frac{k_0}{k_1 + k_0}$$

$\tilde{\theta}$ is the quantile of order $\frac{k_0}{k_1 + k_0}$

Exercise 4. (0-1 loss function leads to the mode of the posterior)

$$\theta \in \Theta \subseteq \mathbb{R}, \forall a \in \Theta, \ell(\theta, a) = \begin{cases} 0 & \theta = a \\ 1 & \theta \neq a \end{cases} : \quad \begin{array}{c} \uparrow \\ \text{---} \\ \bullet \end{array} \quad \begin{array}{c} \rightarrow \\ 1 \\ a \end{array}$$

We're penalizing everything in the same way except one point

We separate two cases: discrete, continuous.

- Discrete case:

$$\begin{aligned} \arg \min_{a \in \Theta} E_{\pi}[\ell(\theta, a) | \underline{x}] &= \arg \min_{a \in \Theta} \sum_{a \in \Theta} \mathbb{1}_{\{\theta \neq a\}} \pi(\theta | \underline{x}) \\ &\stackrel{!}{=} \arg \min_{a \in \Theta} \left(1 - \sum_{a \in \Theta} \mathbb{1}_{\{\theta = a\}} \pi(\theta | \underline{x}) \right) \\ &\stackrel{!}{=} \arg \min_{a \in \Theta} (1 - \pi(a | \underline{x})) \\ &\stackrel{!}{=} \arg \max_{a \in \Theta} \pi(a | \underline{x}) \end{aligned}$$

the point estimate with the 0-1 loss function will be the posterior mode (the mode of the posterior distribution) = point where the posterior has more mass

- Continuous case:

$$\ell_{\varepsilon}(\theta, a) = \begin{cases} 0 & |\theta - a| < \varepsilon \\ 1 & |\theta - a| \geq \varepsilon \end{cases}$$

$$\begin{aligned} \arg \min_{a \in \Theta} E_{\pi}[\ell(\theta, a) | \underline{x}] &= \arg \min_{a \in \Theta} 1 - \int_{\Theta} \mathbb{1}_{(a-\varepsilon, a+\varepsilon)}(\theta) \pi(\theta | \underline{x}) d\theta \\ &\stackrel{!}{=} \arg \max_{a \in \Theta} \int_{\Theta} \mathbb{1}_{(a-\varepsilon, a+\varepsilon)}(\theta) \pi(\theta | \underline{x}) d\theta \\ &\stackrel{!}{=} \arg \max_{a \in \Theta} \pi(da | \underline{x}) \quad \text{for } \varepsilon \rightarrow 0 \end{aligned}$$

→ corresponds again to the posterior mode (mode of the posterior distribution)

Remark: the multinomial distribution is conjugate to the dirichlet distribution

Consider a set of observations: N_1, \dots, N_k ; $n := \sum_{i=1}^k N_i$
(coming from a k -dimensional multinomial distribution)

We set at likelihood: $(N_1, \dots, N_k) \sim \text{multinomial}(\lambda_1, \dots, \lambda_k)$.
We set at prior for $\lambda_1, \dots, \lambda_k$:

$$(\lambda_1, \dots, \lambda_k) \sim \text{dirichlet}(\alpha_1, \dots, \alpha_k)$$

likelihood:

$$p(N_1, \dots, N_k | \lambda_1, \dots, \lambda_k) = \frac{n!}{N_1! \dots N_k!} \lambda_1^{N_1} \dots \lambda_k^{N_k}$$

Prior:

$$\pi(\lambda_1, \dots, \lambda_k) = \frac{1}{B(\underline{\alpha})} \left[\prod_{j=1}^k \lambda_j^{\alpha_j - 1} \right] \mathbb{1}_{\Delta_{k-1}}(\underline{\lambda})$$

We want to show that $\underline{\lambda} | \underline{N} \sim \text{Dirichlet}$:

$$\begin{aligned} \pi(\underline{\lambda} | \underline{N}) &= \frac{p(\underline{N} | \underline{\lambda}) \pi(\underline{\lambda})}{m(\underline{N})} \propto p(\underline{N} | \underline{\lambda}) \pi(\underline{\lambda}) = \\ &= \underbrace{\left[\prod_{j=1}^k \lambda_j^{N_j} \right] \left[\prod_{j=1}^k \lambda_j^{\alpha_j - 1} \mathbb{1}_{\Delta_{k-1}}(\underline{\lambda}) \right]}_{\text{kernel of a dirichlet distribution}} \\ &= \underbrace{\prod_{j=1}^k \lambda_j^{\alpha_j + N_j - 1} \mathbb{1}_{\Delta_{k-1}}(\underline{\lambda})}_{\sim \text{dirichlet}(\alpha_1 + N_1, \dots, \alpha_k + N_k)} \end{aligned}$$

Exercise.

We have an urn full of balls of different colors.

We have $k=3$ different colors. We do a prior guess of the composition:

- 0.1 red balls
- 0.4 blue balls
- 0.5 green balls

We draw $n=100$ balls with replacement and we get:

- $N_1 = 33$ red balls
- $N_2 = 15$ blue balls
- $N_3 = 52$ green balls

$$\Rightarrow \begin{cases} \pi_1(\underline{\lambda}) \stackrel{d}{=} \text{Dirichlet}(0.1, 0.4, 0.5) \\ \pi_1(\underline{\lambda} | \underline{N}) \stackrel{d}{=} \text{Dirichlet}(0.1 + 33, 0.4 + 15, 0.5 + 52) = \text{Dir}(\alpha_1^*, \alpha_2^*, \alpha_3^*) \end{cases}$$

We want to study a bit this distribution:

\Rightarrow point estimate under the quadratic loss function: = posterior mean

$$E_{\pi_1}[\lambda_i | \underline{N}] = \frac{\alpha_i^*}{\alpha_1^* + \alpha_2^* + \alpha_3^*} \Rightarrow \begin{cases} \hat{\lambda}_1 = 0.328 \\ \hat{\lambda}_2 = 0.152 \\ \hat{\lambda}_3 = 0.520 \end{cases} \quad \left. \begin{array}{l} \text{posterior point} \\ \text{estimate for the} \\ \text{composition of} \\ \text{our urn} \end{array} \right.$$

Another prior guess:

- 10 red balls
- 40 blue balls
- 50 green balls

This prior is not only the composition,
it also assumes the number of balls in the urn!

$$\Rightarrow \begin{cases} \pi_2(\underline{\lambda}) = \text{Dir}(10, 40, 50) \\ \pi_2(\underline{\lambda} | \underline{N}) = \text{Dir}(10+33, 40+15, 50+52) \end{cases}$$

Point estimate: $E_{\pi_2}[\lambda_j | \underline{N}] = \frac{\alpha_j^*}{\alpha_1^* + \alpha_2^* + \alpha_3^*}$

$$\Rightarrow \begin{cases} \hat{\lambda}_1 = 0.215 \\ \hat{\lambda}_2 = 0.275 \\ \hat{\lambda}_3 = 0.510 \end{cases} \neq \text{previous ones}$$

Exercise 1. NORMAL-INV-GAMMA MODEL (continuous data)

8/10/2020

1/3

$$X_1, \dots, X_n | \mu, \sigma^2 \sim N(\mu, \sigma^2)$$

$$(\mu, \sigma^2) \sim \text{Normal-inv-gamma}(m_0, k_0, a_0, b_0)$$

$$\mu | \sigma^2 \sim N(m_0, \frac{\sigma^2}{k_0}) \stackrel{d}{=} (2\pi)^{-\frac{1}{2}} \left(\frac{\sigma^2}{k_0} \right)^{-\frac{1}{2}} e^{-\frac{k_0}{2\sigma^2} (\mu - m_0)^2} \mathbb{1}_{(-\infty, \infty)}(\mu)$$

$$\sigma^2 \sim \text{inv-gamma}(a_0, b_0) \stackrel{d}{=} \frac{b_0^{a_0}}{\Gamma(a_0)} (\sigma^2)^{-a_0-1} e^{-\frac{b_0}{\sigma^2}} \mathbb{1}_{(0, \infty)}(\sigma^2)$$

A normal-inverse-gamma, conditionally to a sample that it Gaussian distributed is still a normal-inverse-gamma.

Likelihood:

$$L(\mu, \sigma^2; \underline{x}) = \prod_{i=1}^n f_{X_i}(\underline{x}_i, \mu, \sigma^2) = [\dots]$$

Prior:

$$\pi(\mu, \sigma^2)$$

Posterior:

$$\pi(\mu, \sigma^2 | \underline{x}) \propto \pi(\mu, \sigma^2) L(\mu, \sigma^2, \underline{x}) \\ \sim \text{normal-inverse-gamma}(m_n, k_n, a_n, b_n)$$

$$\begin{cases} m_n = \frac{k_0 m_0 + n \bar{x}}{k_0 + n} \\ k_n = k_0 + n \\ a_n = a_0 + \frac{n}{2} \\ b_n = b_0 + \frac{1}{2} \sum_{i=1}^n (x_i - \bar{x})^2 + \frac{1}{2} \cdot \frac{k_0 n}{k_0 + n} (\bar{x} - m_0)^2 \end{cases}$$

How to choose m_0, k_0 or a_0, b_0 ?

One way is to rely on the summaries of the posterior distributions and see what changes when we modify the parameters of the priors.

For example: let's elicit the parameters of the distribution part of μ .

Elicit the distribution of μ :

$$\begin{cases} \mu | \sigma^2 \sim N(m_0, \frac{\sigma^2}{k_0}) & \text{(prior)} \\ \mu | \sigma^2, \underline{x} \sim N(m_n, \frac{\sigma^2}{k_n}) & \text{(posterior)} \end{cases}$$

Taking for instance the post. mean:

$$\Rightarrow E[\mu | \sigma^2, \underline{x}] = m_n = \frac{k_0 m_0 + n \bar{x}}{k_0 + n} = \frac{k_0}{k_0 + n} m_0 + \frac{n}{k_0 + n} \bar{x}$$

prior guess
(completely driven by prior distribution)

maximum likelihood estimator
(completely driven by data)

We see that the parameter k_0 takes the same role as n : we can think about an hypothetical synthetic sample (that plays the role of the prior sample) and the observed sample. The hypothetical synthetic is of size k_0 and with "maximum likelihood estimator" m_0 . This is how we translate the prior.

Elicit the distribution of σ^2 :

$$\begin{cases} \sigma^2 \sim \text{inv-gamma}(a_0, b_0) \\ \sigma^2 | \underline{x} \sim \text{inv-gamma}(a_n, b_n) \end{cases}$$

$$\Rightarrow E[\sigma^2 | \bar{X}] = \frac{b_n}{a_n - 1} = \frac{1}{a_0 + \frac{n}{2} - 1} \left[b_0 + \frac{1}{2} S_n^2 + \frac{1}{2} \frac{k_0 n}{k_0 + n} (\bar{x} - m_0)^2 \right]$$

$S_n^2 = \sum_{i=1}^n (x_i - \bar{x})^2$

$$= \frac{\frac{1}{a_0 + \frac{n}{2} - 1} b_0 + \frac{n-1}{a_0 + \frac{n}{2} - 1} S^2 + \frac{1}{2} \frac{k_0 n}{k_0 + n} (\bar{x} - m_0)^2 \cdot \frac{1}{a_0 + \frac{n}{2} - 1}}{a_0 + \frac{n}{2} - 1}$$

we discard this part ($n \rightarrow \infty$ this is vanishing)

$$= \frac{a_0 - 1}{a_0 + \frac{n}{2} - 1} \left(\frac{b_0}{a_0 - 1} \right) + \frac{\frac{n-1}{2} S^2}{a_0 + \frac{n}{2} - 1}$$

expectation of the prior distribution on the variance term sample variance

] $a_0 - 1$ plays the same role as $(n-1)/2$, and $b_0/(a_0 - 1)$ plays the same role as S^2 (the rest of the interpretation is as before)

Exercise 2: BETA - BERNoulli MODEL (0-1 data)

Elections in Florida.

They interviewed $n = 509$ people, $r = 279$ said they'll vote for Bush.

We model the situation as a Bernoulli model: success = vote for Bush.

$$X_1, \dots, X_n | \theta \sim \text{Be}(\theta) \quad \theta = \text{IP(voting for Bush)}$$

The probability function in this case is:

$$p_X(x|\theta) = \theta^x (1-\theta)^{1-x} \quad x \in \{0,1\}$$

likelihood:

$$L(\theta; \underline{x}) = \prod_{i=1}^n p_X(x_i|\theta) = \theta^r (1-\theta)^{n-r}$$

$$\Rightarrow \hat{\theta}_{MLE} = \arg \max_{\theta \in (0,1)} L(\theta, \underline{x}) = \frac{r}{n} = 0.548 \quad \rightarrow \text{according to this Bush is going to win}$$

What if we include some priors? Always look at the likelihood and try to find something of a similar form (maybe we'll obtain a conjugate prior!)

$$\theta \sim \text{Beta}(a, b) : \pi(\theta) = \frac{1}{B(a,b)} \theta^{a-1} (1-\theta)^{b-1}$$

Posterior:

$$\pi(\theta | \underline{x}) \propto \pi(\theta) L(\theta, \underline{x}) = \underbrace{\theta^{a+r-1} (1-\theta)^{b+n-r-1}}_{\text{kernel of a Beta}(a+r, b+n-r)}$$

where $r = \sum_{i=1}^n x_i$

1. What kind of model are we assuming for the data? (Discarding θ)
= What is the marginal distribution of the data?

$$m(\underline{x}) = \int_{\mathbb{R}} L(\theta, \underline{x}) \pi(d\theta) = \int_{\mathbb{R}} \theta^r (1-\theta)^{n-r} \frac{1}{B(a,b)} \theta^{a-1} (1-\theta)^{b-1} \mathbb{1}_{(0,1)}(\theta) d\theta$$

$$= \frac{B(a_n, b_n)}{B(a, b)} \int_{\mathbb{R}} \frac{1}{B(a_n, b_n)} \theta^{a_n-1} (1-\theta)^{b_n-1} \mathbb{1}_{(0,1)}(\theta) d\theta = 1$$

$\begin{cases} a_n = a_0 + r \\ b_n = b_0 + n - r \end{cases}$

2. Choose a_0 and b_0 .

We look at the posterior distr: $\theta | \underline{x} \sim \text{Beta}(a_n, b_n)$

$$\begin{aligned} \mathbb{E}[\theta | \underline{x}] &= \frac{a_n}{a_n + b_n} = \frac{a_0 + r}{a_0 + b_0 + n} = \frac{a_0}{a_0 + b_0 + n} + \frac{r}{a_0 + b_0 + n} \\ &= \frac{a_0 + b_0}{a_0 + b_0 + n} + \frac{a_0}{a_0 + b_0} + \frac{n}{a_0 + b_0 + n} \quad \frac{r}{n} \\ &\quad \text{---} \quad \text{---} \quad \text{---} \\ &\quad \mathbb{E}[\theta] \text{ (prior)} \quad \hat{\theta}_{\text{MUE}} \text{ (empirical)} \end{aligned}$$

To synthesize:

	prior	empirical
sample size	$a_0 + b_0$	n
#success	a_0	r

Suppose that we have some previous informations:
(a previous round of interviews)

- Percentage of preference for Bush: 0.491
- Standard deviation: 0.022

} let's use the informations of a previous round of interviews for the elicitation ↗

$$\Rightarrow \left\{ \begin{array}{l} \mathbb{E}[\theta] = \frac{a_0}{a_0 + b_0} = 0.491 \\ \text{Var}(\theta) = \frac{a_0 b_0}{(a_0 + b_0)^2 (a_0 + b_0 + 1)} = 0.022^2 \end{array} \right.$$

$$\Rightarrow [..] \Rightarrow \left\{ \begin{array}{l} a_0 \approx 253 = r_{\text{old}} \\ b_0 \approx 262 = n_{\text{old}} \end{array} \right.$$

$$\Rightarrow \theta | \underline{x} \sim \text{Beta}(253 + 279, 262 + (504 - 279)) \sim \text{Beta}(532, 492)$$



We can compute summaries of $\pi(\theta | \underline{x})$:

For example: how much mass is on the part of the support $> \frac{1}{2}$?

$$\mathbb{P}(\theta > \frac{1}{2} | \underline{x}) = \int_{0.5}^1 \pi(d\theta | \underline{x}) = 1 - \text{CDF}_{0.5}(\text{Beta}(532, 492)) \approx 0.8947$$

We started from a prior distribution that was centered in a part of the support, but in force of the empirical evidence of the new sample we moved in another part.

3. Predict the outcome for $n+1$.

$$\begin{aligned} p(x_{n+1} | x_1, \dots, x_n) &= \frac{p(x_1, \dots, x_n, x_{n+1})}{p(x_1, \dots, x_n)} = \frac{\text{Beta}(a_n + x_{n+1}, b_n + 1 - x_{n+1})}{\text{Beta}(a_n, b_n)} \\ &= \frac{\Gamma(a_n + x_{n+1}) \Gamma(b_n + 1 - x_{n+1})}{\Gamma(a_n + b_n + 1)} \cdot \frac{\Gamma(a_n + b_n)}{\Gamma(a_n) \Gamma(b_n)} \end{aligned}$$

Suppose we want to predict if the next person is voting for Bush:

$$\mathbb{P}(x_{n+1} = 1 | \underline{x}) = \frac{\Gamma(a_n + 1) \Gamma(b_n)}{\Gamma(a_n + b_n + 1)} \cdot \frac{\Gamma(a_n + b_n)}{\Gamma(a_n) \Gamma(b_n)}$$

$$P(X_{n+1} = 1 | \underline{x}) = \frac{a_n \Gamma(a_n) \Gamma(a_n + b_n)}{(a_n + b_n) \Gamma(a_n + b_n) \Gamma(a_n)} = \frac{a_n}{a_n + b_n} = E[\theta | \underline{x}] = 0.5159$$

What if we have to predict the outcome for more than 1 person?

$$\begin{aligned}
 p(x_{n+1}, \dots, x_{n+m} | x_1, \dots, x_n) &= \frac{p(x_1, \dots, x_n, x_{n+1}, \dots, x_{n+m})}{p(x_1, \dots, x_n)} \\
 &= \frac{\text{Beta}(a_n + r_m, b_n + (m - r_m))}{\text{Beta}(a_n, b_n)} \quad r_m = \sum_{i=n+1}^m x_i \\
 &= \frac{\Gamma(a_n + r_m) \Gamma(b_n + (m - r_m))}{\Gamma(a_n + b_n + m)} \cdot \frac{\Gamma(a_n + b_n)}{\Gamma(a_n) \Gamma(b_n)} \\
 &= \frac{(a_n + r_m - 1) \dots (a_n) \Gamma(a_n) (b_n + (m - r_m) - 1) \dots (b_n) \Gamma(b_n)}{(a_n + b_n + m - 1) \dots (a_n + b_n) \Gamma(a_n + b_n)} \cdot \frac{\Gamma(a_n + b_n)}{\Gamma(a_n) \Gamma(b_n)} \\
 &\stackrel{(a)_b}{=} \frac{(a_n + b_n)^m}{(a_n + b_n)^m} = \text{predictive distribution of } r_m \text{ successes and } m \text{ further sampling conditioned on the first } n \text{ results.}
 \end{aligned}$$

Exercise 3. POISSON-GAMMA MODEL (counting data)

We suppose to have data: X_1, \dots, X_n , where $X_i \in \{0, 1, 2, \dots\}$

$$X_1, \dots, X_n | \lambda \stackrel{iid}{\sim} p(\lambda);$$

$$X_i | \lambda \sim p(\lambda) \triangleq \frac{e^{-\lambda} \lambda^{x_i}}{x_i!} \quad \lambda \in (0, \infty) : \begin{cases} E[X_i] = \lambda \\ \text{Var}(X_i) = \lambda \end{cases}$$

likelihood:

$$L(\lambda | \underline{x}) = \prod_{i=1}^n p(x_i, \lambda) = \frac{e^{-n\lambda} \lambda^{\sum_{i=1}^n x_i}}{\prod_{i=1}^n x_i!} \propto e^{-n\lambda} \lambda^{S_n}, \quad S_n = \sum_{i=1}^n x_i$$

Prior for λ ?

looking at the likelihood we choose a gamma distribution.

$$\lambda \sim \text{Gamma}(a_0, b_0)$$

$$\pi(\lambda) \triangleq \frac{b_0^{a_0}}{\Gamma(a_0)} \lambda^{a_0-1} e^{-b_0 \lambda} \mathbb{1}_{(0, \infty)}(\lambda) \quad \left(E[\lambda] = \frac{a_0}{b_0}, \text{Var}(\lambda) = \frac{a_0}{b_0^2} \right)$$

In this way the posterior is:

$$\begin{aligned}
 \pi(\lambda | \underline{x}) &\propto \pi(\lambda) L(\lambda | \underline{x}) = \lambda^{a_0-1} e^{-b_0 \lambda} \lambda^{S_n} e^{-n\lambda} \\
 &\stackrel{\perp}{=} \lambda^{(a_0+S_n)-1} e^{-(b_0+n)\lambda} \sim \text{Gamma}(a_0 + S_n, b_0 + n) \\
 &\sim \text{Gamma}(a_n, b_n)
 \end{aligned}$$

1. Elicitation of the parameters of the prior?

$$E[\lambda | \underline{x}] = \frac{a_n}{b_n} = \frac{a_0 + S_n}{b_0 + n} = \frac{b_0}{b_0 + n} \frac{a_0}{b_0} + \frac{n}{b_0 + n} \frac{S_n}{n}$$

$E[\lambda]$ (prior) λ MLE (empirical)

b_0 is playing the role of the sample size of a prior (hypothetical) sample while a_0 plays the role of the sum of all the counts

2. Marginal distribution? (of the data)

8/10/2020

3/3

$$m(x) = \frac{p(x|\lambda) \pi(\lambda)}{\pi(\lambda|x)} = \frac{\left(\frac{1}{x!} e^{-\lambda} \lambda^x\right) \left(\frac{b_0^{a_0}}{\Gamma(a_0)} \lambda^{a_0-1} e^{-b_0 \lambda}\right)}{\frac{b_n^{a_n}}{\Gamma(a_n)} \lambda^{a_n-1} e^{-b_n \lambda}}$$

$$= \frac{1}{x!} \frac{b_0^{a_0}}{b_n^{a_n}} \frac{\Gamma(a_n)}{\Gamma(a_0)} = \binom{a_0 + x - 1}{x} \left(\frac{b_0}{b_0 + 1}\right)^{a_0} \left(\frac{1}{b_0 + 1}\right)^x \Rightarrow X \sim \text{Neg-Bin}(a, b)$$

Marginal distribution
for just one observation

3. Prediction: x_{n+1} ?

$$\begin{aligned} p(X_{n+1}|\underline{x}) &= \int_{\Lambda} p(X_{n+1}|\lambda) \pi(\lambda|\underline{x}) d\lambda \\ &= \int_{\Lambda} \frac{1}{X_{n+1}!} e^{-\lambda} \lambda^{X_{n+1}} \frac{b_n^{a_n}}{\Gamma(a_n)} \lambda^{a_n-1} e^{-b_n \lambda} d\lambda \\ &= \frac{1}{X_{n+1}!} \frac{b_n^{a_n}}{\Gamma(a_n)} \frac{\pi(a_n + X_{n+1})}{(b_n + 1)^{a_n + X_{n+1}}} \int_{\Lambda} [\dots] d\lambda \\ &= \binom{a_n + X_{n+1} - 1}{X_{n+1}} \left(\frac{b_n}{b_n + 1}\right)^{a_n} \left(\frac{1}{b_n + 1}\right)^{X_{n+1}} \end{aligned}$$

\Rightarrow The predictive distribution of the $n+1$ observation given x_1, \dots, x_n is still a negative binomial distribution

$$\sim \text{Neg-Bin}(a_n, b_n)$$

$$X \sim \text{Neg-Bin}(a, b) \stackrel{d}{=} \binom{a+x-1}{x} \left(\frac{b}{b+1}\right)^a \left(\frac{1}{b+1}\right)^x$$

HYPOTHESIS TESTING

$$H_0 \text{ vs. } H_1 : \begin{cases} H_0: \theta \in \Theta_0 \\ H_1: \theta \in \Theta_1 \end{cases} \quad \Theta_0 \cap \Theta_1 = \emptyset \quad \Theta_0 \cup \Theta_1 = \Theta = \text{parameters space}$$

In frequentist statistic we assume the null hypothesis to be true and we look at the empirical evidence and we state if we can reject or not H_0 . The test is not exchangeable w.r.t. the hypothesis. In bayesian statistic the framework is different and we can switch the hypothesis: we're not assuming any of the hypothesis true, we just test the empirical evidence in favor of one hypothesis or the other. The test is equivalent if we switch the hypothesis.

We define a prior:

$$\pi(\theta) = \underbrace{\pi_0 g_0(\theta) \mathbb{1}_{\Theta_0}}_{\text{prior weights that we give to the sets } \Theta_0 \text{ and } \Theta_1} + \underbrace{(1-\pi_0) g_1(\theta) \mathbb{1}_{\Theta_1}}$$

$g_0(\theta), g_1(\theta)$ are the restrictions of the prior to Θ_0, Θ_1

What we can do is compute:

- Prior odds $o_1 := \frac{\pi_0}{1-\pi_0}$
- Posterior odds $o_2 := \frac{P(\Theta_0 | X)}{P(\Theta_1 | X)}$
- Bayesian Factor: $BF_{01} := \frac{\text{posterior odds } o_2}{\text{prior odds } o_1}$

This quantity measures how much the empirical evidence is shifting the mass in favor of the null hypothesis discarding the prior guess

BF_{01}	$z \log BF_{01}$	evidence
<1	<0	in favor of H_1
1-3	0-2	weak evidence for H_0
3-12	2-5	in favor of H_0
12-150	5-10	strong in favor of H_0
>150	>10	very strong in favor of H_0

How to compute the Bayes factor:

$$BF_{01} = \frac{\frac{P(\Theta_0 | X)}{P(\Theta_1 | X)}}{\frac{P(X | \Theta_0)}{P(X | \Theta_1)}} \cdot \frac{\frac{P(\Theta_1)}{P(\Theta_0)}}{\frac{P(X | \Theta_1)}{P(X | \Theta_0)}} = \frac{P(\Theta_0, X)}{P(X)} \cdot \frac{P(X)}{P(\Theta_1, X)} \cdot \frac{P(\Theta_1)}{P(\Theta_0)}$$

we can define the Bayes factor also in the case in which Θ_0 is not a set but a single point

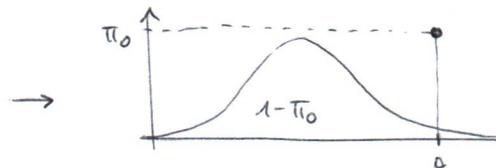
$$H_0: \theta = \theta_0$$

$$H_1: \theta \neq \theta_0$$

$$\begin{cases} H_0: \theta = \theta_0 \\ H_1: \theta \neq \theta_0 \end{cases}$$

We define a prior:

$$\pi_1(\theta) = \pi_0 \delta_{\theta_0} + (1 - \pi_0) g(\theta) \mathbb{1}_{\{\theta \neq \theta_0\}}$$



Within this framework:

$$P(H_0) = P(\theta = \theta_0) = \pi_0$$

$$P(H_1) = P(\theta \neq \theta_0) = 1 - \pi_0$$

Marginal distribution under this prior:

$$\begin{aligned} m_1(x) &= w_1(x_1, \dots, x_n) = \int_{\mathbb{R}^n} f(x|\theta) \pi_1(\theta) d\theta \\ &\stackrel{!}{=} \int_{\mathbb{R}^n} f(x|\theta) \pi_0 \delta_{\theta_0} d\theta + \int_{\mathbb{R}^n} f(x|\theta) \pi_0 \mathbb{1}_{\{\theta \neq \theta_0\}} d\theta \\ &= f(x|\theta_0) \pi_0 + (1 - \pi_0) \underline{m(x)} \end{aligned}$$

marginal of X under
the general prior $\pi(\theta)$

$$P(\theta = \theta_0 | x) = \frac{f(x|\theta_0) \pi_0}{f(x|\theta_0) \pi_0 + (1 - \pi_0) m(x)}$$

Now we can compute:

$$\begin{aligned} BF_{01} &= \frac{\text{posterior odds}_{01}}{\text{prior odds}_{01}} = \frac{P(\theta = \theta_0 | x)}{P(\theta \neq \theta_0 | x)} \cdot \frac{1 - \pi_0}{\pi_0} \\ &\stackrel{!}{=} \frac{\frac{f(x|\theta_0) \pi_0}{m_1(x)}}{1 - \frac{f(x|\theta_0) \pi_0}{m_1(x)}} \cdot \frac{1 - \pi_0}{\pi_0} \\ &\stackrel{!}{=} \frac{\frac{f(x|\theta_0) \pi_0}{m_1(x)}}{\frac{m_1(x) - f(x|\theta_0) \pi_0}{m_1(x)}} \cdot \frac{1 - \pi_0}{\pi_0} \\ &\stackrel{!}{=} \frac{(1 - \pi_0) f(x|\theta_0)}{f(x|\theta_0) \pi_0 + (1 - \pi_0) m(x) - f(x|\theta_0) \pi_0} \\ &\stackrel{!}{=} \frac{f(x|\theta_0)}{m(x)} \end{aligned}$$

with the def.
of $m_1(x)$

Doing Hypothesis Testing through Bayes Factor is equivalent to a model comparison:

$$\text{Model 0: } \begin{cases} X_1, \dots, X_n | \theta \sim \pi_0(\theta) \\ \pi_0(\theta) \stackrel{d}{=} \delta_{\theta_0} \end{cases}$$

Here we compare a single mass point against a set, we can use it more generally: (with a restricted set)

$$H_0: \theta = \theta_0$$

$$H_1: \theta < \theta_0$$

$$\text{Model 1: } \begin{cases} X_1, \dots, X_n | \theta \sim \pi_1(\theta) \\ \pi_1(\theta) \stackrel{d}{=} \pi(\theta) \leftarrow \text{starting prior that we use} \end{cases}$$

$$\pi^*(\theta) \propto \pi(\theta) \mathbb{1}_{(\theta < \theta_0)}$$

Exercise 1.

Back to Beta-Bernoulli model: (for the elections (Bush))

$$\begin{cases} X_1, \dots, X_n | \theta \sim \text{Be}(\theta) \\ \pi(\theta) \sim \text{Beta}(a, b) \end{cases}$$

We found out:

$$\text{IP}(\theta > 0.5 | \underline{x}) = 0.8928$$

Suppose we want to test: $\begin{cases} H_0: \theta \leq 0.5 \\ H_1: \theta > 0.5 \end{cases}$

We compute the Bayes factor:

$$BF_{01} = \frac{\text{IP}(\theta \leq 0.5 | \underline{x})}{\text{IP}(\theta > 0.5 | \underline{x})} \cdot \frac{\text{IP}(\theta > 0.5)}{\text{IP}(\theta \leq 0.5)} = \frac{1 - 0.8928}{0.8928} \cdot \frac{1 - 0.6543}{0.6543} = 0.0634$$

$$2 \log(BF_{01}) = -5.517 \implies \text{empirical evidence in favor of } H_1$$

Exercise 2.

We flip a coin $n=1000$ times, obtaining 452 heads.

Question: is the coin balanced?

$$\begin{cases} H_0: p = 0.5 \\ H_1: p \neq 0.5 \end{cases} \quad p = \text{IP}(\text{success})$$

We assume a binomial distribution.

$$X = \# \text{ successes} \\ X \sim \text{Bin}(n, p) \stackrel{d}{=} \binom{n}{x} p^x (1-p)^{n-x}$$

We start from:

$$\text{p}(x | p_0) = \binom{n}{x} p_0^x (1-p_0)^{n-x} \stackrel{H_0}{=} \binom{n}{x} p_0^n \quad (\text{since } H_0: p_0 = 1 - p_0 = \frac{1}{2})$$

We need the marginal for the data:

we first need to assume a prior for p :

$$\pi(p) \sim \text{Beta}(1, 1) \sim U([0, 1])$$

with this assumption:

$$m(x) = \binom{n}{x} \frac{\text{Beta}(1+x, 1+(n-x))}{\text{Beta}(1, 1)}$$

Now we can compute:

$$BF_{01} = \frac{\text{p}(x | \theta_0)}{m(x)} = \frac{\binom{n}{x} p_0^n \text{Beta}(1, 1)}{\binom{n}{x} \text{Beta}(1+x, 1+n-x)} \\ = p_0^n \frac{\binom{n}{x} (n+1)!}{n!} = 0.2512$$

$$2 \log(BF_{01}) = -2.763 < 0 \implies \text{evidence for } H_1 \\ \implies \text{The coin is } \underline{\text{not}} \text{ balanced}$$

Exercise 3.

$$f(x|\theta) = 2\theta e^{-\theta x^2} \mathbb{1}_{(0,+\infty)}(x)$$

$$X_1, \dots, X_n | \theta \stackrel{iid}{\sim} f(x, \theta)$$

1. Find the family of prior conjugate to $f(x|\theta)$ and compute the posterior.

We start from the likelihood:

$$L(\theta, \underline{x}) = \prod_{i=1}^n f(x_i, \theta) = 2^n \theta^n \left(\prod_{i=1}^n x_i \right) e^{-\theta \sum_{i=1}^n x_i^2} := S_{n,2}$$

we're interested
only in θ

$\underbrace{\theta^n e^{-\theta S_{n,2}}}_{\text{it reminds a kernel
of a gamma distribution}}$

$$\text{We try with: } \theta \sim \text{Gamma}(a, b)$$

$$\begin{aligned} \pi(\theta | \underline{x}) &\propto \pi(\theta) L(\theta, \underline{x}) = \left[\theta^{a-1} e^{-b\theta} \right] \left[\theta^n e^{-\theta S_{n,2}} \right] \\ &= \theta^{a+n-1} e^{-(b+S_{n,2})\theta} \sim \text{Gamma}(a+n, b+S_{n,2}) \end{aligned}$$

2. Write the Bayes Factor to test: $H_0: \theta = 1, H_1: \theta \neq 1$

It's a pivotal (single valued) test.

Marginal of the data:

$$\begin{aligned} m(\underline{x}) &= \int_{\Theta} L(\theta, \underline{x}) \pi(\theta) d\theta \\ &= \int_{\Theta} 2^n \theta^n \left(\prod_{i=1}^n x_i \right) e^{-\theta S_{n,2}} \frac{b^a}{\Gamma(a)} \theta^{a-1} e^{-b\theta} d\theta \\ &= 2^n \left(\prod_{i=1}^n x_i \right) \frac{b^a}{\Gamma(a)} \frac{\Gamma(a+n)}{(b+S_{n,2})^{a+n}} \int_{\Theta} \cdot d\theta \end{aligned}$$

We can now compute:

$$\begin{aligned} BF_{01} &= \frac{f(\underline{x} | \theta_0)}{m(\underline{x})} = \frac{2^n \theta_0^n \left(\prod_{i=1}^n x_i \right) e^{-\theta_0 S_{n,2}}}{2^n \left(\prod_{i=1}^n x_i \right) \frac{b^a}{\Gamma(a)} \frac{\Gamma(a+n)}{(b+S_{n,2})^{a+n}}} \\ &= \frac{\theta_0^n e^{-\theta_0 S_{n,2}} \Gamma(a) (b+S_{n,2})^{a+n}}{b^a \Gamma(a+n)} \\ &= \frac{e^{-S_{n,2}} \Gamma(a) (b+S_{n,2})^{a+n}}{b^a \Gamma(a+n)} \quad \theta_0 = 1 \end{aligned}$$

3. Suppose that we have a sample of $n=4$, with: $S_{4,2} = 5.71$

a. Choose a, b s.t. $E[\theta] = 1, \text{Var}[\theta] = 10$

b. Test $H_0: \theta_0 = 1, H_1: \theta_0 \neq 1$

a. $\begin{cases} E[\theta] = \frac{a}{b} = 1 \\ \text{Var}[\theta] = \frac{a}{b^2} = 10 \end{cases} \Rightarrow \begin{cases} a = 0.1 \\ b = 0.1 \end{cases} \Rightarrow \pi(\theta) \stackrel{d}{=} \text{Gamma}(0.1, 0.1)$

b. $\begin{aligned} BF_{01} &= \frac{e^{-5.71} \Gamma(0.1)}{0.1^{0.1} \Gamma(0.1+4)} (0.1+5.71)^{0.1+4} \\ &= 7.91 \quad \Rightarrow \text{in favor of } H_0 \text{ (since } > 1) \end{aligned}$

4. Compute: $P(X_5 > 1 | \underline{x})$.

9/10/2020

3/3

First we derive the CDF of X :

$$\begin{aligned} F(x) &= \int_{-\infty}^x 2\theta + e^{-\theta t^2} \mathbb{1}_{(0,+\infty)}(t) dt \\ &= [\dots] = 1 - e^{-\theta x^2} \end{aligned}$$

$$\begin{aligned} P(X_5 > 1 | \underline{x}) &= \int_{\underline{x}}^{\infty} P(X_5 > 1 | \theta) \pi(\theta | \underline{x}) d\theta \\ &= \int_0^{\infty} (1 - F(1)) \pi(\theta | \underline{x}) d\theta \\ &= \int_0^{\infty} e^{-\theta} \frac{(b + s_{n,2})^{a+n}}{\Gamma(a+n)} \theta^{a+n-1} e^{-(b+s_{n,2})\theta} d\theta \\ &\stackrel{?}{=} \frac{(b + s_{n,2})^{a+n}}{(b + s_{n,2} + 1)^{a+n}} \int_{\underline{x}}^{\infty} \dots d\theta \\ &= \left(\frac{b + s_{n,2}}{b + s_{n,2} + 1} \right)^{a+n} = 0.5215 \end{aligned}$$

JEFFREY'S PRIOR - Binomial case

= prior proportional to the square root of the Fisher information matrix

We suppose to have:

$$X \sim \text{Bin}(n, p) : p(x|p) = \binom{n}{x} p^x (1-p)^{n-x}$$

Fisher information matrix:

$$I(\underline{\theta}) = \mathbb{E}\left[\left(\frac{\partial}{\partial \underline{\theta}} \log f(x, \underline{\theta})\right)^2 \middle| \underline{\theta}\right] = -\mathbb{E}\left[\frac{\partial^2}{\partial \underline{\theta}^2} \log f(x, \underline{\theta}) \middle| \underline{\theta}\right]$$

So:

$$\log p(x|p) \propto x \log(p) + (n-x) \log(1-p)$$

$$\frac{d}{dp} \log p(x|p) = \frac{x}{p} - \frac{n-x}{1-p}$$

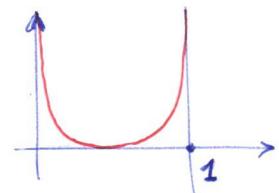
$$\frac{d^2}{dp^2} \log p(x|p) = -\frac{x}{p^2} - \frac{n-x}{(1-p)^2}$$

$$\begin{aligned} I(p) &= -\mathbb{E}\left[-\frac{x}{p^2} - \frac{n-x}{(1-p)^2}\right] = \frac{np}{p^2} + \frac{n-np}{(1-p)^2} \\ &= \frac{n-np+np}{p(1-p)} = \frac{n}{p(1-p)} \end{aligned}$$

\Rightarrow Jeffreys prior:

$$\pi_j(p) \propto \sqrt{I(p)} \propto \sqrt{\frac{1}{p(1-p)}} = p^{-1/2} (1-p)^{-1/2}$$

$\underbrace{\quad}_{\text{Beta}(\frac{1}{2}, \frac{1}{2})}$



JEFFREY'S PRIOR

family of objective priors
= non informative priors

$$\text{prior } \propto \sqrt{|I(\theta)|}, \quad I(\theta) = \mathbb{E} \left[\left(\frac{\partial}{\partial \theta} \log f(x, \theta) \right)^2 | \theta \right] = \mathbb{E} \left[\left(\frac{f'(x, \theta)}{f(x, \theta)} \right)^2 | \theta \right]$$

$$= - \mathbb{E} \left[\frac{\partial^2}{\partial \theta^2} \log f(x, \theta) | \theta \right]$$

$$\text{Moreover: } I_n(\theta) = n I(\theta)$$

Univariate Normal Distribution
(with prior on both parameters)

$$X_1, \dots, X_n | \mu, \sigma^2 \sim N(\mu, \sigma^2)$$

$$f(x, \mu, \sigma^2) \propto (\sigma^2)^{-1/2} e^{-\frac{1}{2\sigma^2}(x-\mu)^2}$$

$$\log f(x, \mu, \sigma^2) \propto -\frac{1}{2} \log(\sigma^2) - \frac{1}{2\sigma^2}(x-\mu)^2$$

$$\frac{d}{d\mu} \log f(x, \mu, \sigma^2) = \frac{(x-\mu)}{2\sigma^2}$$

$$\frac{d^2}{d\mu^2} \log f(x, \mu, \sigma^2) = -\frac{1}{2\sigma^2}$$

$$\frac{d}{d\sigma^2} \log f(x, \mu, \sigma^2) = -\frac{1}{2\sigma^2} + \frac{(x-\mu)^2}{2(\sigma^2)^2}$$

$$\frac{d^2}{d(\sigma^2)^2} \log f(x, \mu, \sigma^2) = \frac{1}{2(\sigma^2)^2} - \frac{(x-\mu)^2}{(\sigma^2)^3}$$

$$\frac{d^2}{d\mu d\sigma^2} \log f(x, \mu, \sigma^2) = -\frac{(x-\mu)}{2(\sigma^2)^2}$$

$$I(\theta) = \mathbb{E}_X \begin{bmatrix} -\frac{\partial^2}{\partial \mu^2} & -\frac{\partial^2}{\partial \mu \partial \sigma^2} \\ -\frac{\partial^2}{\partial \mu \partial \sigma^2} & -\frac{\partial^2}{\partial (\sigma^2)^2} \end{bmatrix} = \mathbb{E}_X \begin{bmatrix} \frac{1}{2\sigma^2} & \frac{x-\mu}{2(\sigma^2)^2} \\ \frac{x-\mu}{2(\sigma^2)^2} & \frac{1}{2(\sigma^2)^2} - \frac{(x-\mu)^2}{(\sigma^2)^3} \end{bmatrix}$$

$$= \begin{bmatrix} \frac{1}{2\sigma^2} & 0 \\ 0 & \frac{1}{2(\sigma^2)^2} - \frac{\sigma^2}{(\sigma^2)^3} \end{bmatrix} = \begin{bmatrix} \frac{1}{\sigma^2} & 0 \\ 0 & \frac{1}{2(\sigma^2)^2} \end{bmatrix}$$

$$|I(\theta)| = \det I(\theta) = \frac{1}{\sigma^2} \cdot \frac{1}{2(\sigma^2)^2} = \frac{1}{2(\sigma^2)^3}$$

$$\Rightarrow \sqrt{|I(\theta)|} \propto (\sigma^2)^{-3/2}$$

$$\Rightarrow \underbrace{\pi_j(\theta)}_{\text{Jeffreys prior}} \propto (\sigma^2)^{-3/2}$$

Jeffreys prior

proportional on something depending on σ^2
 \Rightarrow for the location parameter it proportional to a constant (when the only parameter we consider is the location parameter))

Moreover, this is an improper prior (it is not normalized integrating up to 1): it is not normalized and cannot be normalized)

Quick remark for the location family case

$$X_1, \dots, X_n | \mu \stackrel{iid}{\sim} f(x, \mu)$$

$$f(x, \mu) \in \mathcal{F}_\mu = \{g(x-\mu) : \mu \in \mathbb{R}, g(\cdot) \text{ density}\}$$

the act of the location parameter
is in term of translation
(we just translate x by μ)

We have that:

$$\frac{d}{d\mu} g(x-\mu) = -g'(x-\mu)$$

Moreover:

$$\frac{f'(x, \mu)}{f(x, \mu)} = -\frac{g'(x-\mu)}{g(x-\mu)}$$

$$\begin{aligned} \Rightarrow I(\theta) &= I(\mu) = \mathbb{E} \left[\left(\frac{f'(\mu)}{f(\mu)} \right)^2 \middle| \mu \right] \\ &= \int \frac{(g'(x-\mu))^2}{(g(x-\mu))^2} g(x-\mu) dx \quad \text{because we're taking the expectation w.r.t. } x \\ &= \int \frac{(g'(z))^2}{g(z)} dz \quad \underbrace{\text{c constant} := k}_{\text{since it does not depend on } \mu} \end{aligned}$$

$$\Rightarrow \Pi_j(\mu) \propto \text{constant}$$

Jeffrey's prior

Multivariate Normal Distribution

Normal likelihood with Normal-Inverse-Wishart prior

Riccardo Corradin

October 14, 2020

1/13

N - NIW

\mathbf{X} is a random vector, absolutely continuous wrt the Lebesgue measure on \mathbb{R}^p , distributed as a multivariate Normal distribution, $\mathbf{X} \sim N(\mu, \Sigma)$, with $\mathbf{X} \in \mathbb{R}^p$, $\mu \in \mathbb{R}^p$ and $\Sigma \in \mathbb{M}_{p \times p}$, with $\mathbb{M}_{p \times p}$ denoting the space of positive definite symmetric matrices $p \times p$. The density function of \mathbf{X} is

$$f(\mathbf{x}; \mu, \Sigma) = (2\pi)^{-d/2} |\Sigma|^{-1/2} \exp \left\{ -\frac{1}{2} (\mathbf{x} - \mu)^\top \Sigma^{-1} (\mathbf{x} - \mu) \right\}$$

with $\mathbb{E}[\mathbf{X}] = \mu$ and $\text{Cov}(X_i, X_j) = \Sigma_{ij}$. We assume as prior distributions $\pi(\mu, \Sigma)$ on (μ, Σ) a Normal-Inverse-Wishart distribution, $(\mu, \Sigma) \sim NIW(\mathbf{m}_0, \kappa_0, \nu_0, \Lambda_0)$, which corresponds to assume that

$$\begin{aligned} \mu | \Sigma &\sim N(\mathbf{m}_0, \Sigma / \kappa_0), & \mathbf{m}_0 \in \mathbb{R}^p, \kappa_0 \in \mathbb{R}^+ \\ \Sigma &\sim IW(\nu_0, \Lambda_0), & \nu_0 \in \mathbb{R}^+, \nu_0 > p - 1, \Lambda_0 \in \mathbb{M}_{p \times p} \end{aligned}$$

2/13

N - NIW

Recall if $\Sigma \sim IW(\nu_0, \Lambda_0)$, with $\Sigma, \Lambda_0 \in \mathbb{M}_{p \times p}$, $\nu_0 \in \mathbb{R}^+$ and $\nu_0 > p - 1$, its density is equal to

$$f(\Sigma; \nu_0, \Lambda_0) = \frac{|\Lambda_0|^{\nu_0/2}}{2^{\nu_0 p/2} \Gamma_p(\nu_0/2)} |\Sigma|^{-(\nu_0 + p + 1)/2} \exp \left\{ -\frac{1}{2} \text{tr}(\Lambda_0 \Sigma^{-1}) \right\}$$

where $\Gamma_p(\cdot)$ denotes the multivariate Gamma function of order p with

$$\Gamma_p(\nu_0/2) = \pi^{p(p-1)/4} \prod_{j=1}^p \Gamma \left(\frac{\nu_0}{2} + \frac{1-j}{2} \right).$$

We further recall that $\mathbb{E}[\Sigma] = \Lambda_0 / (\nu_0 - p - 1)$ and

$$\text{Cor}(\Sigma_{ij}, \Sigma_{kl}) = \frac{2\Lambda_{0[ij]} \Lambda_{0[kl]} + (\nu_0 - p - 1)(\Lambda_{0[ik]} \Lambda_{0[jl]} + \Lambda_{0[il]} \Lambda_{0[jk]})}{(\nu_0 - p)(\nu_0 - p - 1)^2 (\nu_0 - p - 3)}$$

3/13

N - NIW

(Recall:)

Multivariate location-scale t-Student distribution

We further recall that $\mathbf{Y} \sim t_\alpha(\mu, \Sigma)$ denotes a multivariate location-scale t-Student distribution, with $\alpha > 0$ degree of freedom, $\mu \in \mathbb{R}^p$ location parameter and $\Sigma \in \mathbb{M}_{p \times p}$ scale parameter, and its density is equal to

$$f(\mathbf{y}; \alpha, \mu, \Sigma) = \frac{\Gamma_p(\frac{\alpha+p}{2})}{\Gamma(\frac{\alpha}{2}) \alpha^{p/2} \pi^{p/2} |\Sigma|^{1/2}} \left[1 + \frac{1}{\alpha} (\mathbf{y} - \mu)^\top \Sigma^{-1} (\mathbf{y} - \mu) \right]^{-(\alpha-p)/2}$$

4/13

Details (Normal-Inverse-Wishart)

N - NIW

The full model specification is then

$$\begin{aligned} \mathbf{X}_1, \dots, \mathbf{X}_n | \boldsymbol{\mu}, \Sigma &\sim N(\boldsymbol{\mu}, \Sigma) \\ \boldsymbol{\mu} | \Sigma &\sim N(\mathbf{m}_0, \Sigma / \kappa_0) \\ \Sigma &\sim IW(\nu_0, \Lambda_0) \end{aligned}$$

where a priori we have

$$\mathbb{E}_{\pi}[\boldsymbol{\mu}] = \mathbb{E}_{\Sigma}[\mathbb{E}_{\boldsymbol{\mu}}[\boldsymbol{\mu} | \Sigma]] = \mathbb{E}_{\Sigma}[\mathbf{m}_0] = \mathbf{m}_0$$

$$\begin{aligned} \underline{Var}_{\pi}(\boldsymbol{\mu}) &= Var_{\Sigma}(\mathbb{E}_{\boldsymbol{\mu}}[\boldsymbol{\mu} | \Sigma]) + \mathbb{E}_{\Sigma}[Var_{\boldsymbol{\mu}}(\boldsymbol{\mu} | \Sigma)] \\ &= Var_{\Sigma}(\mathbf{m}_0) + \mathbb{E}_{\Sigma}[\Lambda_0 / \kappa_0] = \frac{1}{\kappa_0} \left(\frac{\Lambda_0}{\nu_0 - p - 1} \right) \end{aligned}$$

and

$$\mathbb{E}_{\pi}[\Sigma] = \Lambda_0 / (\nu_0 - p - 1); \quad \underline{Var}_{\pi}(\Sigma) = ..see\ above$$

5/13

N - NIW

The likelihood is equal to

$$\begin{aligned} L(\boldsymbol{\mu}, \Sigma; \mathbf{X}_{1:n}) &= \prod_{i=1}^n f_X(\mathbf{X}_i; \boldsymbol{\mu}, \Sigma) \\ &= \prod_{i=1}^n (2\pi)^{-p/2} |\Sigma|^{-1/2} \exp \left\{ -\frac{1}{2} (\mathbf{x}_i - \boldsymbol{\mu})^\top \Sigma^{-1} (\mathbf{x}_i - \boldsymbol{\mu}) \right\} \\ &= (2\pi)^{-np/2} |\Sigma|^{-n/2} \exp \left\{ -\frac{1}{2} \sum_{i=1}^n (\mathbf{x}_i - \boldsymbol{\mu})^\top \Sigma^{-1} (\mathbf{x}_i - \boldsymbol{\mu}) \right\} \end{aligned}$$

we further set $\mathbf{S}_n = \sum_{i=1}^n \mathbf{x}_i$.

6/13

N - NIW

The posterior distribution of $(\boldsymbol{\mu}, \Sigma)$ given a sample $\mathbf{X}_1, \dots, \mathbf{X}_n$ is equal to

$$\begin{aligned} \pi(\boldsymbol{\mu}, \Sigma | \mathbf{X}_1, \dots, \mathbf{X}_n) &\propto L(\boldsymbol{\mu}, \Sigma; \mathbf{X}_{1:n}) \pi(\boldsymbol{\mu}, \Sigma) \\ &\propto \underbrace{|\Sigma|^{-1/2} |\Sigma|^{-(\nu_0+n+p+1)/2}}_A \exp \left\{ -\frac{1}{2} \sum_{i=1}^n (\mathbf{X}_i - \boldsymbol{\mu})^\top \Sigma^{-1} (\mathbf{X}_i - \boldsymbol{\mu}) \right. \\ &\quad \left. - \frac{1}{2} (\boldsymbol{\mu} - \mathbf{m}_0)^\top \left(\frac{\Sigma}{\kappa_0} \right)^{-1} (\boldsymbol{\mu} - \mathbf{m}_0) - \frac{1}{2} \text{tr}(\Lambda_0 \Sigma^{-1}) \right\} \\ &= A \exp \left\{ -\frac{1}{2} \left[\sum_{i=1}^n \mathbf{X}_i^\top \Sigma^{-1} \mathbf{X}_i - 2\boldsymbol{\mu}^\top \Sigma^{-1} \mathbf{S}_n + n\boldsymbol{\mu}^\top \Sigma^{-1} \boldsymbol{\mu} \right. \right. \\ &\quad \left. \left. + \kappa_0 \boldsymbol{\mu}^\top \Sigma^{-1} \boldsymbol{\mu} - 2\kappa_0 \boldsymbol{\mu}^\top \Sigma^{-1} \mathbf{m}_0 + \kappa_0 \mathbf{m}_0^\top \Sigma^{-1} \mathbf{m}_0 + \text{tr}(\Lambda_0 \Sigma^{-1}) \right] \right\} \end{aligned}$$

7/13

N - NIW

We set $\kappa_n = \kappa_0 + n$ and $\mathbf{m}_n = (\kappa_0 \mathbf{m}_0 + \mathbf{S}_n) / \kappa_n$

$$\begin{aligned} &= A \exp \left\{ -\frac{1}{2} \left[-2\boldsymbol{\mu}^\top \left(\frac{\Sigma}{\kappa_n} \right)^{-1} \mathbf{m}_n + \boldsymbol{\mu}^\top \left(\frac{\Sigma}{\kappa_n} \right)^{-1} \boldsymbol{\mu} \pm \mathbf{m}_n^\top \left(\frac{\Sigma}{\kappa_n} \right)^{-1} \mathbf{m}_n \right. \right. \\ &\quad \left. \left. + \sum_{i=1}^n \mathbf{X}_i^\top \Sigma^{-1} \mathbf{X}_i + \kappa_0 \mathbf{m}_0^\top \Sigma^{-1} \mathbf{m}_0 + \text{tr}(\Lambda_0 \Sigma^{-1}) \right] \right\} \\ &= A \exp \left\{ -\frac{1}{2} \left[(\boldsymbol{\mu} - \mathbf{m}_n)^\top \left(\frac{\Sigma}{\kappa_n} \right)^{-1} (\boldsymbol{\mu} - \mathbf{m}_n) \right. \right. \\ &\quad \left. \left. + \text{tr} \left(\left(\sum_{i=1}^n \mathbf{X}_i \mathbf{X}_i^\top + \kappa_0 \mathbf{m}_0 \mathbf{m}_0^\top - \kappa_n \mathbf{m}_n \mathbf{m}_n^\top + \Lambda_0 \right) \Sigma^{-1} \right) \right] \right\} \end{aligned}$$

8/13

$$\begin{aligned}
&= |\Sigma|^{-1/2} \exp \left\{ -\frac{1}{2} \left[(\mu - \mathbf{m}_n)^\top \left(\frac{\Sigma}{\kappa_n} \right)^{-1} (\mu - \mathbf{m}_n) \right] \right\} \\
&\times |\Sigma|^{-(\nu_0 + p + 1)/2} \exp \left\{ -\frac{1}{2} \left[\text{tr} \left(\left(\Lambda_0 + \sum_{i=1}^n (\mathbf{X}_i - \bar{\mathbf{X}})(\mathbf{X}_i - \bar{\mathbf{X}})^\top + \frac{\kappa_0 \eta}{\kappa_n} (\bar{\mathbf{X}} - \mathbf{m}_0)(\bar{\mathbf{X}} - \mathbf{m}_0)^\top \right) \Sigma^{-1} \right) \right] \right\}
\end{aligned}$$

which is the kernel of a NIW distribution with updated parameters

$$\begin{aligned}
\kappa_n &= \kappa_0 + n, & \mathbf{m}_n &= (\kappa_0 \mathbf{m}_0 + \mathbf{S}_n) / \kappa_n, & \nu_n &= \nu_0 + n, \\
\Lambda_n &= \Lambda_0 + \sum_{i=1}^n (\mathbf{X}_i - \bar{\mathbf{X}})(\mathbf{X}_i - \bar{\mathbf{X}})^\top + \frac{\kappa_0 \eta}{\kappa_n} (\bar{\mathbf{X}} - \mathbf{m}_0)(\bar{\mathbf{X}} - \mathbf{m}_0)^\top
\end{aligned}$$

9/13

N - NIW

The marginal distribution of μ is equal to

$$\begin{aligned}
\pi(\mu) &= \int_{\mathbb{M}_{p \times p}} \pi(\mu | \Sigma) \pi(d\Sigma) \\
&= \int_{\mathbb{M}_{p \times p}} (2\pi)^{-p/2} \left| \frac{\Sigma}{\kappa_0} \right|^{-1/2} \exp \left\{ -\frac{1}{2} (\mu - \mathbf{m}_0)^\top \left(\frac{\Sigma}{\kappa_0} \right)^{-1} (\mu - \mathbf{m}_0) \right\} \\
&\times \frac{|\Lambda_0|^{\nu_0/2}}{2^{(\nu_0 p)/2} \Gamma_p(\frac{\nu_0}{2})} |\Sigma|^{-(\nu_0 + p + 1)/2} \exp \left\{ -\frac{1}{2} \text{tr}(\Lambda_0 \Sigma^{-1}) \right\} d\Sigma \\
&\propto |\Lambda_0|^{\nu_0/2} \int_{\mathbb{M}_{p \times p}} |\Sigma|^{-(\nu_0 + p + 1)/2} \exp \left\{ -\frac{1}{2} \text{tr}[(\kappa_0(\mu - \mathbf{m}_0)(\mu - \mathbf{m}_0)^\top + \Lambda_0)\Sigma^{-1}] \right\} d\Sigma
\end{aligned}$$

10/13

N - NIW

$$\begin{aligned}
&\propto \frac{|\Lambda_0|^{\nu_0/2} 2^{(\nu_0+1)p/2} \Gamma_p(\frac{\nu_0+1}{2})}{|\Lambda_0 + \kappa_0(\mu - \mathbf{m}_0)(\mu - \mathbf{m}_0)^\top|^{(\nu_0+1)/2}} \\
&\propto \left| \Lambda_0 + (\mu - \mathbf{m}_0)^\top \left(\frac{\Lambda_0}{\kappa_0} \right)^{-1} (\mu - \mathbf{m}_0) \Lambda_0 \right|^{-(\nu_0+1)/2} \\
&\propto \left| 1 + \frac{1}{(\nu_0 - p + 1)} (\mu - \mathbf{m}_0)^\top \left(\frac{\Lambda_0}{(\nu_0 - p + 1)\kappa_0} \right)^{-1} (\mu - \mathbf{m}_0) \right|^{-(\nu_0+1)/2}
\end{aligned}$$

which is the kernel of a t -Student distribution, $\mu \sim t_{\nu_0 - p + 1} \left(\mathbf{m}_0, \frac{\Lambda_0}{(\nu_0 - p + 1)\kappa_0} \right)$.

11/13

N - NIW

Given $\mathbf{X}_{1:n} = \{\mathbf{X}_1, \dots, \mathbf{X}_n\}$ a set of n observations, the predictive distribution of the $n+1$ observation is equal to

$$\begin{aligned}
f(\mathbf{X}_{n+1} | \mathbf{X}_{1:n}) &= \int_{\mathbb{R}^p \times \mathbb{M}_{p \times p}} f(\mathbf{X}_{n+1} | \mu, \Sigma, \mathbf{X}_{1:n}) \pi(d\mu, d\Sigma) \\
&= \int_{\mathbb{R}^p \times \mathbb{M}_{p \times p}} f(\mathbf{X}_{n+1} | \mu, \Sigma) \pi(d\mu, d\Sigma | \mathbf{X}_{1:n}) \\
&\propto \int_{\mathbb{R}^p \times \mathbb{M}_{p \times p}} |\Sigma|^{-1/2} \exp \left\{ -\frac{1}{2} (\mathbf{x}_{n+1} - \mu)^\top \Sigma^{-1} (\mathbf{x}_{n+1} - \mu) \right\} \\
&\times \left| \frac{\Sigma}{\kappa_n} \right|^{-1/2} \exp \left\{ -\frac{1}{2} (\mu - \mathbf{m}_n)^\top \left(\frac{\Sigma}{\kappa_n} \right)^{-1} (\mu - \mathbf{m}_n) \right\} \\
&\times |\Sigma|^{-(\nu_n + p + 1)/2} \exp \left\{ -\frac{1}{2} \text{tr}(\Lambda_n \Sigma^{-1}) \right\} d\Sigma d\mu
\end{aligned}$$

12/13

$$\propto C \int_{\mathbb{R}^p \times \mathbb{M}_{p \times p}} C^{-1} |\Sigma|^{-1/2} \exp \left\{ -\frac{1}{2} (\mu - \mathbf{m}_{n+1})^\top \left(\frac{\Sigma}{\kappa_{n+1}} \right)^{-1} (\mu - \mathbf{m}_{n+1}) \right\} \\ \times |\Sigma|^{-(\nu_{n+1} + p + 1)/2} \exp \left\{ -\frac{1}{2} \text{tr}(\Lambda_{n+1} \Sigma^{-1}) \right\} d\Sigma d\mu$$

which is the integral of the kernel of a NIW, up to a normalization constant C , the part of C depending on \mathbf{X}_{n+1} is

$$C \propto \left| \Lambda_n + \frac{\kappa_n}{\kappa_n + 1} (\mathbf{x}_{n+1} - \mu_n)(\mathbf{x}_{n+1} - \mu_n)^\top \right|^{-(\nu_{n+1}/2)} \\ \propto \left| 1 + \frac{1}{\nu_n - p + 1} (\mathbf{x}_{n+1} - \mu_n)^\top \left(\frac{(\kappa_n + 1)\Lambda_n}{\kappa_n(\nu_n - p + 1)} \right)^{-1} (\mathbf{x}_{n+1} - \mu_n) \right|^{-(\nu_{n+1}/2)}$$

that describes the kernel of a multivariate t-Student distribution,

$$\mathbf{X}_{n+1} | \mathbf{X}_1, \dots, \mathbf{X}_n \sim t_{\nu_n - p + 1} \left(\mathbf{m}_n, \left(\frac{(\kappa_n + 1)\Lambda_n}{\kappa_n(\nu_n - p + 1)} \right) \right).$$

13/13

Recap:

- Normal-Inverse-Wishart prior
- Normal likelihood
- Normal-Inverse-Wishart posterior
- The marginal on μ is a t-student (multivariate)
- The predictive is a t-student (multivariate)

Exercise 1.

15/10/2020

4/4

3 components working in parallel ($\Rightarrow \Pi$).

They're of the same kind. We denote by Y_1, Y_2, Y_3 the failure time of each component (expressed in hours) and we assume : $Y_i \sim \text{Beta}(a, 1)$, $a > 0$ $\forall i$.

We collected a sample of $n = 120$ and we obtained $\sum_{i=1}^n \log(x_i) = -4a$, where $x_i = \text{lifetime of the } i\text{-th system}$.

- Derive the density of X (= lifetime of the entire system).

$$F_X(x) = P(X \leq x) = P(\max\{Y_1, Y_2, Y_3\} \leq x) = P(Y_1 \leq x, Y_2 \leq x, Y_3 \leq x) \\ = P(Y_1 \leq x) P(Y_2 \leq x) P(Y_3 \leq x) = (P(Y \leq x))^3$$

$$Y_i \sim \text{Beta}(a, 1) \implies f(y, a) = a y^{a-1} \underset{(0,1)}{\mathbb{1}}(y) \\ \implies F_Y(y) = \int_0^y a t^{a-1} dt = y^a$$

$$\implies F_X(x) = (P(Y \leq x))^3 = (x^a)^3 = x^{3a}$$

$$\implies f(x, a) = \frac{d}{dx} F_X(x) = 3a x^{3a-1} \implies x | a \sim \text{Beta}(3a, 1)$$

- Find the conjugate distribution and the update rule of the hyperparameters.

$$L(a, x) = \prod_{i=1}^n f(x_i, a) = \prod_{i=1}^n (3a x_i^{3a-1}) = 3^n a^n \left(\prod_{i=1}^n x_i^{-1} \right) e^{3a \sum_{i=1}^n \log(x_i)} \\ \propto a^n e^{3a \sum_{i=1}^n \log(x_i)} = a^n e^{3a S_{\text{ex}}}$$

$$\implies \pi(a) \stackrel{d}{=} \text{Gamma}(\alpha, \beta) \stackrel{d}{=} \frac{\beta^\alpha}{\Gamma(\alpha)} a^{\alpha-1} e^{-\beta a}$$

$$\implies \pi(a|x) \propto L(a, x) \pi(a) \propto a^{\alpha+n-1} e^{-(\beta-3S_{\text{ex}})a} \sim \text{Gamma}(\alpha+n, \beta-3S_{\text{ex}})$$

- Suppose that we have a previous study with $m = 10$ observations :

$$z_1, \dots, z_{10} \text{ st. } \sum_{i=1}^m z_i = -1.73.$$

Use the equivalence sample principle to specify α and β .

We start from the posterior expectation:

$$E[a|x] = \left(\frac{\beta-3S_{\text{ex}}}{\alpha+n} \right)^{-1} = \left(\frac{\alpha}{\alpha+n} + \frac{n}{\alpha+n} \left(-\frac{3S_{\text{ex}}}{n} \right) \right)^{-1}$$

so n is playing the role of α ,

$-3S_{\text{ex}}$ is playing the role of β .

A way to elicitate the hyperparameters is :

$$\alpha = m = 10$$

$$\beta = -3S_{\text{ex}} = 5.19$$

The parameters of the posterior distribution are :

$$\alpha_n = \alpha + n = m + n = 10 + 120 = 130$$

$$\beta_n = \beta - 3S_{\text{ex}} = 5.19 - \dots = 154.95$$

parameters of the posterior distr.
assuming as prior elicitation the
specification according to the
equivalence sample principle
with the prior sample z_1, \dots, z_m

- Calculate the predictive probability that a system works at most 30 minutes.

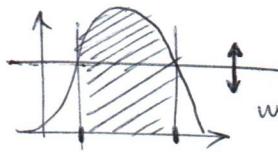
$$P(X_{n+1} \leq 0.5|x) = \int_0^{\infty} P(X_{n+1} \leq 0.5|a) \pi(a|x) da \\ = \int_0^{\infty} (0.5)^{3a} a^{\alpha+n-1} e^{-(\beta-3S_{\text{ex}})a} \frac{\beta_n^{\alpha_n}}{\Gamma(\alpha_n)} da$$

$$F_x(0.5)$$

$$P(X_{n+1} \leq 0.5 | \underline{x}) = \frac{\beta_n^{\alpha_n}}{\Gamma(\alpha_n)} \cdot \frac{\Gamma(\alpha_n)}{(\beta_n - 3\log(0.5))^{\alpha_n}} \int_0^{\infty} \frac{(\beta_n - 3\log(0.5))^{\alpha_n}}{\Gamma(\alpha_n)} a^{\alpha_n-1} e^{-(\beta_n - 3\log(0.5))a} da$$

$$= \frac{\beta_n^{\alpha_n}}{(\beta_n - 3\log(0.5))^{\alpha_n}} = 0.1768$$

5. Use an opportune approximation to derive a 0.95 HPD credible interval.



we have a distribution, we move the line until we have a mass of 0.95 inside the interval

For large values of α , by CLT:

$\text{Egamma} \xrightarrow{d} N(\text{E}[a], \text{Var}(a))$

$$\text{In this case: } E[a|\underline{x}] = \frac{\alpha_n}{\beta_n} = \frac{130}{154} = 0.8390$$

$$\text{Var}(a|\underline{x}) = \frac{\alpha_n}{\beta_n^2} = (0.0736)^2$$

$$\Rightarrow CI_{0.95}(a) = E[a|\underline{x}] \pm z_{0.025} \sqrt{\text{Var}(a|\underline{x})} = (0.7028, 0.9772)$$

6. Test $H_0: \alpha \geq 1$ vs. $H_1: \alpha < 1$ using the Bayes Factor.

[Hint: CDF $\chi^2_{20}(10.38) = 0.0393$]

$$BF_{01} = \frac{\text{posterior odds}_{01}}{\text{prior odds}_{01}} = \frac{P(\alpha \geq 1 | \underline{x})}{P(\alpha < 1 | \underline{x})} \cdot \frac{P(\alpha < 1)}{P(\alpha \geq 1)}$$

A prior: $a \sim \text{Egamma}(\alpha, \beta) = \text{Egamma}(10, 5.19) = \text{Egamma}\left(\frac{20}{2}, 5.19\right)$

And so: $c \cdot a \sim \text{Egamma}\left(\frac{20}{2}, \frac{1}{2}\right)$ if $c = 2\beta = 2 \cdot 5.19 = 10.38$

$(c \cdot a \sim \text{Egamma}(a, \frac{\beta}{c}))$

$$\Rightarrow P(\alpha < 1) = P\left(\text{Egamma}\left(\frac{20}{2}, \frac{1}{2}\right) < 10.38\right) = P\left(\chi^2_{20} < 10.38\right) = 0.0393$$

for the posterior:

$$P(\alpha < 1 | \underline{x}) \approx \phi\left(\frac{1 - 0.8390}{0.0736}\right) = 0.9856$$

$$\Rightarrow BF_{01} = \frac{1 - 0.9856}{0.9856} \cdot \frac{0.0393}{1 - 0.0393} = 0.0005 \dots$$

\Rightarrow strong evidence in support for H_1

$\Rightarrow \alpha < 1$

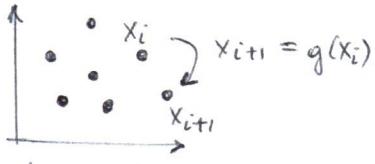
RANDOM NUMBERS

We're dealing with pseudo random numbers.

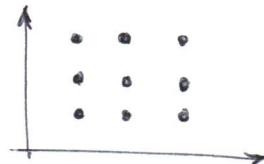
How do we obtain random numbers?

- Sample in some way a sequence from a uniform: $U_n \sim U([0, 1])$
- Something cool
- Deal with other distributions

There is a distinction: pseudo random vs. quasi random



pseudo random: x_{i+1} is generated deterministically



quasi random: equally distributed forming a grid

RNG (Random Numbers Generator)

Defined by a structure: (S, S_0, f, U, g) :

- S = support = finite set of state
- S_0 = initial condition
- f = transition kernel: $f: S \rightarrow S$
- U = set of outputs
- g = output function: $g: S \rightarrow U$

Pseudocode:

1. Initialize S_0
2. for $i = 1, \dots, n_{\text{rep}}$:
 - $S_i = f(S_{i-1})$
 - $U_i = g(S_i)$

Period of a random number generator: smallest m s.t. $S_{i+m} = S_i \quad \forall i$

Example: Linear Congruential Generators

$$X_n = (a X_{n-1} + c) \bmod m$$

Exercise 1.

Suppose we have X_1, \dots, X_n taking values in \mathbb{R}^+ .

$$X_1, \dots, X_n | \theta, \beta \sim \text{Weibull}(\theta, \beta)$$

$$f(x; \theta, \beta) = \frac{\beta}{\theta} x^{\beta-1} \exp\left\{-\frac{x^\beta}{\theta}\right\} \mathbb{1}_{\mathbb{R}^+}(x)$$

Trick for the conjugate prior (general variable θ):

$$L(\theta, \cdot) \propto \theta^n e^{-\theta} \Rightarrow \text{gamma}$$

$$\propto \frac{1}{\theta^n} e^{-\theta} \Rightarrow \text{IG}$$

inverse gamma

1. Suppose β known. Derive the conjugate prior $\Pi(\theta)$ and write the update rule for the hyperparameters.

$$L(\theta; x, \beta) = \prod_{i=1}^n f(x_i; \theta, \beta) = \frac{\beta^n}{\theta^n} \left[\prod_{i=1}^n x_i^{\beta-1} \right] e^{-\sum_{i=1}^n \frac{x_i^\beta}{\theta}} \mathbb{1}_{\mathbb{R}^+}(x_i)$$

$$= \frac{\beta^n}{\theta^n} \left[\prod_{i=1}^n x_i^{-1} \right] e^{-\sum_{i=1}^n \frac{x_i^\beta}{\theta}} + \beta \sum_{i=1}^n \log(x_i) \mathbb{1}_{\mathbb{R}^+}(x_i)$$

$$\propto \frac{1}{\theta^n} e^{-\frac{1}{\theta} \sum_{i=1}^n x_i^\beta} := S(\beta, x)$$

$$\propto \frac{1}{\theta^n} e^{-\frac{1}{\theta} S(\beta, x)} \quad \Rightarrow \text{looks like the kernel of an inverse gamma distribution, we set:}$$

$$\theta \sim \text{IG}(a, b)$$

$\theta \sim IG(a, b)$:

$$\pi(\theta) = \frac{b^a}{\Gamma(a)} \theta^{-a-1} e^{-\frac{b}{\theta}}$$

$$\pi(\theta | \mathbf{x}) \propto L(\theta; \mathbf{x}, \beta) \pi(\theta) \propto \theta^{-a-n-1} e^{-\frac{1}{\theta}(b + S(\beta, \mathbf{x}))}$$

$$\Rightarrow \theta | \mathbf{x} \sim IG(a_n, b_n) \quad a_n = a+n, \quad b_n = b + S(\beta, \mathbf{x})$$

$$\Rightarrow \begin{cases} X | \beta, \theta & \text{Weibull} \\ \theta | \beta & \text{Inverse Gamma } (a, b) \\ \theta | X, \beta & \text{Inverse Gamma } (a_n, b_n) \end{cases}$$

2. Suppose we have a sample collected from a previous study: $\mathbf{z} = (z_1, \dots, z_m)$. Specify the hyperparameters of $\pi(\theta)$ using the Equivalent Sample Principle (ESP). We start from the MLE: (for θ)

$$\log f(\theta; \mathbf{x}, \beta) \propto -n \log(\theta) - \frac{S(\beta, \mathbf{x})}{\theta}$$

$$\Rightarrow \frac{\partial}{\partial \theta} \log L(\theta; \mathbf{x}, \beta) = -\frac{n}{\theta} + \frac{S(\beta, \mathbf{x})}{\theta^2} = 0 \Rightarrow \hat{\theta}_{MLE} = \frac{S(\beta, \mathbf{x})}{n}$$

Moreover we have:

$$E[\theta | \mathbf{x}] = \frac{bn}{a+n-1} = \frac{b + S(\beta, \mathbf{x})}{a+n-1} = \frac{a-1}{a+n-1} \cdot \frac{b}{a-1} + \frac{n}{a+n-1} \cdot \frac{S(\beta, \mathbf{x})}{n}$$

$a-1$ plays the role of the sample size
 b plays the role of $S(\beta, \mathbf{x})$

$$\Rightarrow \text{Elicitation : } \begin{cases} a-1 = m \\ b = S(\beta, \mathbf{x}) \end{cases} \Rightarrow \begin{cases} a_1 = m+1 \\ b_1 = S(\beta, \mathbf{x}) \end{cases}$$

- 2'. Choose the hyperparameters of the conjugate prior assuming $a_2 > 0$ fixed and $b_2 = \arg \max_b m(\mathbf{z}; a_2, b)$

We need the marginal distribution:

$$\begin{aligned} m(\mathbf{z}; a_2, b) &= \int L(\theta; \mathbf{z}, \beta) \pi(\theta) d\theta \\ &\stackrel{?}{=} \int \frac{b^m}{\theta^m} \left(\prod_{i=1}^m z_i^{\beta-1} \right) e^{-\frac{S(\beta, \mathbf{z})}{\theta}} \frac{b^{a_2}}{\Gamma(a_2)} \theta^{-a_2-1} e^{-\frac{b}{\theta}} d\theta \\ &= \left(\prod_{i=1}^m z_i^{\beta-1} \right) b^m \frac{b^{a_2}}{\Gamma(a_2)} \cdot \frac{\Gamma(a_2+m)}{(b+S(\beta, \mathbf{z}))^{a_2+m}} \int \theta^{-a_2-m-1} e^{-\frac{1}{\theta}(b+S(\beta, \mathbf{z}))} \cdot \frac{(b+S(\beta, \mathbf{z}))^{a_2+m}}{\Gamma(a_2+m)} d\theta \end{aligned}$$

$$\frac{d}{db} m(\mathbf{z}; a_2, b) \propto \frac{d}{db} \left(\frac{b^{a_2}}{(b+S(\beta, \mathbf{z}))^{a_2+m}} \right) = \frac{a_2 b^{a_2-1} (b+S(\beta, \mathbf{z}))^{a_2+m} - b^{a_2} (a_2+m) (b+S(\beta, \mathbf{z}))^{a_2+m-1}}{(b+S(\beta, \mathbf{z}))^{2(a_2+m)}}$$

$$\frac{d}{db} \left(\frac{b^{a_2}}{(b+S(\beta, \mathbf{z}))^{a_2+m}} \right) = 0 \Rightarrow b_2 = \frac{a_2 S(\beta, \mathbf{z})}{m}$$

3. Use the Bayes factor to compare, based on the sample \mathbf{x} , two models with the same likelihood and:

$$M_1 : \pi_1(\theta) = \pi(\theta; a_1, b_1) \quad a_1 = m+1, \quad b_1 = S(\beta, \mathbf{z})$$

$$M_2 : \pi_2(\theta) = \pi(\theta; a_2, b_2) \quad a_2 = m, \quad b_2 = \frac{a_2 S(\beta, \mathbf{z})}{m}$$

The idea is to test which specification of the hyperparameter is better.

$$\begin{aligned}
 BF_{12} &= \frac{m(\underline{x}, a_1, b_1)}{m(\underline{x}, a_2, b_2)} = \frac{\beta^n (\prod_{i=1}^n x_i^{\beta-1})}{\beta^n (\prod_{i=1}^n x_i^{\beta-1})} \cdot \frac{\frac{\Gamma(a_1+n)}{\Gamma(a_1)} \cdot \frac{b_1^{a_1}}{(b_1 + S(\beta, \underline{x}))^{a_1+n}}}{\frac{\Gamma(a_2+n)}{\Gamma(a_2)} \cdot \frac{b_2^{a_2}}{(b_2 + S(\beta, \underline{x}))^{a_2+n}}} \\
 &= \frac{\Gamma(m+n+1) (S(\beta, \underline{x}))^{m+n}}{\Gamma(m+n) (S(\beta, \underline{x}) + S(\beta, \underline{x}))^{m+n+1}} \cdot \frac{\Gamma(m) (S(\beta, \underline{x}) + S(\beta, \underline{x}))^{m+n}}{\Gamma(m+n) (S(\beta, \underline{x}))^m} \\
 &= \frac{m+n}{m} \cdot \frac{S(\beta, \underline{x})}{S(\beta, \underline{x}) + S(\beta, \underline{x})}
 \end{aligned}$$

Suppose: $\underline{z} = (1.51, 0.50, 1.20, 0.32)$

$\underline{x} = (0.54, 0.69, 5.47, 5.29, 4.13, 9.51, 3.68, 4.96, 3.70, 1.54)$

$\beta = 2$

and so: $S(\beta, \underline{z}) = 4.0725$
 $S(\beta, \underline{x}) = 150.2753$

$$\Rightarrow BF_{12} = \frac{4+10}{4} \cdot \frac{4.0725}{4.0725 + 150.2753} = 0.0923$$

$$\Rightarrow z \log(BF_{12}) = -4.7654$$

} we prefer the second prior

Exercise 2.

We consider a M/M/1 model.

$X_i = T_i - T_{i-1}$ = interarrival time :

Y_i = service time :

$$X_1, \dots, X_n | \lambda \sim \mathcal{E}(\lambda) \quad \mathbb{E}[X] = \frac{1}{\lambda}$$

$$Y_1, \dots, Y_n | \mu \sim \mathcal{E}\left(\frac{1}{\mu}\right)$$

Data consists in pairs: $(X_1, Y_1), \dots, (X_n, Y_n)$.

Crucial information: $(\underline{X}) \perp\!\!\!\perp (\underline{Y})$.

1. Define the likelihood and derive the conjugate prior $\pi(\lambda, \mu)$. Moreover find the hyperparameters of the posterior.

$$\begin{aligned}
 L(\lambda, \mu | \underline{X}, \underline{Y}) &= \prod_{i=1}^n f(x_i, y_i; \lambda, \mu) = \prod_{i=1}^n f(x_i; \lambda) f(y_i; \mu) \\
 &= \prod_{i=1}^n \lambda \mu e^{-\lambda x_i} e^{-\mu y_i} = \lambda^n \mu^n e^{-\sum \lambda x_i} e^{-\sum \mu y_i} \\
 &= \lambda^n \mu^n e^{-\lambda \sum x_i} e^{-\mu \sum y_i}
 \end{aligned}$$

$$\pi(\lambda, \mu) = \pi(\mu) \pi(\lambda) :$$

$$\pi(\lambda) \stackrel{def}{=} \frac{b^a}{\Gamma(a)} \lambda^{a-1} e^{-b\lambda}$$

$$\lambda \sim \text{Gamma}(a, b), \quad \mu \sim \text{Gamma}(c, d)$$

$$\pi(\mu) \stackrel{def}{=} \frac{d^c}{\Gamma(c)} \mu^{c-1} e^{-d\mu}$$

$$\text{Posterior: } \lambda | \underline{X}, \underline{Y} \sim \lambda | \underline{X} \sim \text{Gamma}(a+n, b+S_x)$$

$$\mu | \underline{X}, \underline{Y} \sim \mu | \underline{Y} \sim \text{Gamma}(c+n, d+S_y)$$

- 1! Elicitation of the prior parameters considering an old sample ($n=8$):

$$\underline{x}_{\text{old}} = (2, 1.8, 1.7, 2.5, 2.1, 2.0, 1.9, 2.2)$$

$$\underline{y}_{\text{old}} = (1.7, 1.5, 0.9, 1.4, 1.8, 1.6, 1.0, 0.8)$$

$$\mathbb{E}[\lambda | \underline{x}] = \frac{\alpha n}{b n} = \frac{\alpha + n}{b + S_x} = \left(\frac{b + S_x}{\alpha + n} \right)^{-1} = \left(\frac{\alpha}{\alpha + n} \frac{b}{a} + \frac{n}{\alpha + n} \frac{S_x}{n} \right)^{-1}$$

a plays the same role as n
 b plays the same role as S_x
 we have the same with $\mathbb{E}[\mu | \underline{y}]$

Elicitation:

$$\begin{aligned} a &= c = n \\ b &= S_{x \text{ old}} \\ d &= S_{y \text{ old}} \end{aligned}$$

2. Test the hypothesis: $H_0 : \lambda = 0.5, \mu = 1, H_1 : \lambda \neq 0.5, \mu \neq 1$

$$BF_{01} = \frac{\prod_{i=1}^n f(x_i, y_i | \lambda=0.5, \mu=1)}{m(\underline{x}, \underline{y})} \quad \leftarrow \text{likelihood with } H_0$$

$$\begin{aligned} m(\underline{x}, \underline{y}) &= \iint \prod_{i=1}^n f(x_i, y_i | \lambda, \mu) \pi(\lambda, \mu) d\mu d\lambda \\ &= \frac{b^a d^c}{\Gamma(a) \Gamma(c)} \iint \lambda^n \mu^n \lambda^{a-1} \mu^{c-1} e^{-(b+S_x)\lambda - (d+S_y)\mu} d\mu d\lambda \\ &= \frac{b^a d^c}{\Gamma(a) \Gamma(c)} \frac{\Gamma(a+n) \Gamma(c+n)}{(b+S_x)^{a+n} (d+S_y)^{c+n}} \cancel{\iint \dots d\mu d\lambda} \end{aligned}$$

Considering: $n = 12, S_x = 22.6, S_y = 18.4$

$$\Rightarrow BF_{01} = \frac{L(\lambda=0.5, \mu=1 | \underline{x}, \underline{y})}{m(\underline{x}, \underline{y})} = 0.4363 \Rightarrow \text{since it's close to 1 there is no big evidence in case of } H_0 \text{ or } H_1$$

3. We have that a system is stable if $\frac{\lambda}{\mu} < 1$.

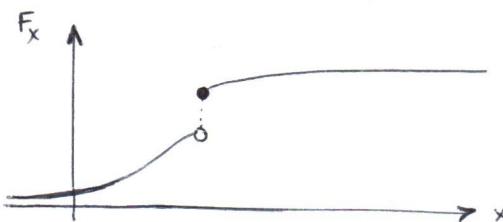
Compute the posterior expectation of $\frac{\lambda}{\mu}$.

$$\mathbb{E}\left[\frac{\lambda}{\mu} | \underline{x}, \underline{y}\right] = \underbrace{\mathbb{E}[\lambda | \underline{x}]}_{\text{gamma}} \cdot \underbrace{\mathbb{E}\left[\frac{1}{\mu} | \underline{y}\right]}_{\text{inverse gamma}} = \left(\frac{\alpha + n}{b + S_x}\right) \left(\frac{d + S_y}{c + n - 1}\right)$$

INVERSE TRANSFORM METHOD

We want to sample $X \sim F_X(\cdot)$ with $F_X(\cdot)$ = cumulative distr. function:

$F: \mathbb{R} \rightarrow [0, 1]$ monotonically increasing, right continuous, $F(-\infty) = 0$, $F(+\infty) = 1$



(we can have discontinuities)

$$F(x) \sim U([0, 1])$$

We consider the generalized inverse of $F(x)$: $U = F^{-1}(X) = \inf \{x : F(x) \geq U\}$

Prop. If $U \sim U([0, 1]) \Rightarrow X = F^{-1}(U) \sim F_X$

Proof.

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

The proposition says: $A_1 = A_2$.

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

$$u \leq F(x_0) \Rightarrow u \in A_2 \Rightarrow A_1 \subseteq A_2$$

Reversely, consider $u \in A_2$: $u \leq F(x_0)$. By the def. of F : $F^{-1}(u) \leq x_0$
 $\Rightarrow u \in A_1 \Rightarrow A_2 \subseteq A_1 \Rightarrow A_1 = A_2$.

$$\begin{aligned} P(X \leq x_0) &= P(F^{-1}(U) \leq x_0) &= P(U \leq F(x_0)) = F(x_0) \\ &\stackrel{\text{construction}}{\uparrow} & \uparrow & \uparrow \\ A_1 = A_2 & & U \sim U([0, 1]) & \blacksquare \end{aligned}$$

Algorithm:

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

1. Sample $U_i \sim U([0, 1])$
2. Set $X_i = F^{-1}(U_i)$

Example 1. (Exponential distribution)

$X \sim \mathcal{E}(\lambda)$; $f(x) = \lambda e^{-\lambda x} \mathbb{1}_{(0, \infty)}(x)$ continuous w.r.t. the Lebesgue measure
 We want to sample from this distribution.

$$F(x) = \int_0^x \lambda e^{-\lambda t} \mathbb{1}_{(0, \infty)}(t) dt = 1 - e^{-\lambda x}$$

We impose:

$$\begin{aligned} U = F(X) &= 1 - e^{-\lambda X} & \Rightarrow 1 - U = e^{-\lambda X} \\ & \Rightarrow \log(1 - U) = -\lambda X & \Rightarrow X = -\frac{1}{\lambda} \log(1 - U) \stackrel{d}{=} -\frac{1}{\lambda} \log(U) \end{aligned}$$

$$\begin{aligned} U &\sim U([0, 1]) = \\ 1 - U &\sim U([0, 1]) \end{aligned}$$

Example 2. (Gamma distribution)

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

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

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

Algorithm:

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

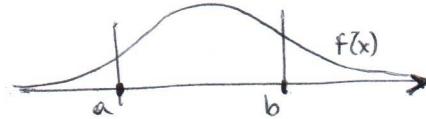
1. Sample $Y_1, \dots, Y_m \sim \text{Gamma}(1, b) = \mathcal{E}(b)$ (Example 1.)
2. Set $X_i = \sum_{j=1}^m Y_j$

Example 3. (Truncated distribution)

Truncated distribution.

We have $X \sim F_x(\cdot)$ absolutely continuous.

We want to sample from $Y = X |_{X \in A}$, where $A = (a, b)$.



$$F_Y(y) = F_{X|_{X \in A}}(y) = P(X \leq y | X \in A) = \frac{P(X \leq y, X \in A)}{P(X \in A)}$$

$$= \begin{cases} \frac{F_x(y) - F_x(a)}{F_x(b) - F_x(a)} & a < y \leq b \\ 0 & y \leq a \\ 1 & y > b \end{cases}$$

What is the density?

$$f_Y(y) = \frac{dF_Y}{dy} = \frac{1}{P(Y \in A)} f_X(y) \mathbb{1}_{(a,b)}(x)$$

How to sample?

$$U = F_Y(y) = \frac{F_x(y) - F_x(a)}{F_x(b) - F_x(a)} := \frac{F_x(y) - F_1}{F_2 - F_1} \rightarrow \underbrace{(F_2 - F_1)U + F_1}_{\sim U([F_1, F_2])} = F_x(y)$$

\Rightarrow instead of sampling from $U([0,1])$, for a truncated version we sample from $U([F(a), F(b)])$

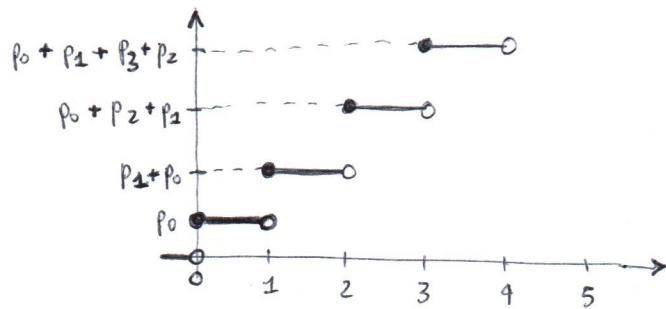
Example 4. (Discrete distribution)

We want to sample from a discrete random variable.

$$X \sim F_x(\cdot), X \in \{0, 1, 2, \dots, n\}, n \in \mathbb{N}$$

$$\text{and to: } p_X(x) = \sum_{j=0}^n p_j \delta_j(x).$$

Graphically:



$$F_X(x) = \sum_{j=0}^n p_j \mathbb{1}_{\{j \leq x\}}$$

$$F_X^{-1}(U) = \begin{cases} 0 & U < p_0 \\ 1 & p_0 \leq U < p_0 + p_1 \\ 2 & p_0 + p_1 \leq U < p_0 + p_1 + p_2 \\ \vdots & \vdots \\ n & \sum_{j=0}^{n-1} p_j \leq U \leq \sum_{j=0}^n p_j \end{cases}$$

Algorithm:

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

1. Sample $U \sim U([0,1])$
2. Proceed according to \circledast

However, we cannot use this method to sample from a gaussian.

BOX-MULLER METHOD

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

We start from:

$$(X_1, X_2) \stackrel{iid}{\sim} N(0,1) : f(x_1, x_2) = (2\pi)^{-1/2} e^{-\frac{1}{2}(x_1^2 + x_2^2)}$$

We move to polar coordinates:

$$\begin{cases} \rho = x_1^2 + x_2^2 & \in [0, \infty) \\ \varphi = \tan^{-1}\left(\frac{x_1}{x_2}\right) & \in [0, 2\pi) \end{cases} \iff \begin{cases} x_1 = \sqrt{\rho} \cos(\varphi) \\ x_2 = \sqrt{\rho} \sin(\varphi) \end{cases}$$

$$J = \begin{bmatrix} \frac{\partial x_1}{\partial \rho} & \frac{\partial x_1}{\partial \varphi} \\ \frac{\partial x_2}{\partial \rho} & \frac{\partial x_2}{\partial \varphi} \end{bmatrix} = \begin{bmatrix} \frac{1}{2\sqrt{\rho}} \cos(\varphi) & -\frac{1}{2\sqrt{\rho}} \sin(\varphi) \\ -\frac{1}{2\sqrt{\rho}} \sin(\varphi) & \frac{1}{2\sqrt{\rho}} \cos(\varphi) \end{bmatrix}$$

$$\begin{aligned} \Rightarrow f(\rho, \varphi) &= f_{\underline{X}}(\sqrt{\rho} \cos(\varphi), \sqrt{\rho} \sin(\varphi)) |J| \mathbb{1}_{[0, \infty)}(\rho) \mathbb{1}_{[0, 2\pi)}(\varphi) \\ &\stackrel{!}{=} (2\pi)^{-1/2} e^{-\frac{1}{2}\rho} |J| \\ &\stackrel{!}{=} (2\pi)^{-1/2} e^{-\frac{1}{2}\rho} \frac{1}{2} \\ &\sim U([0, 2\pi]) \quad \sim \Sigma\left(\frac{1}{2}\right) \end{aligned}$$

Algorithm:

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

1. Sample $U_1, U_2 \sim U([0, 1])$
2. Set $\rho = -2\log(U_1)$, $\varphi = (2\pi)U_2$
3. Set $X_1 = \sqrt{\rho} \cos(\varphi)$, $X_2 = \sqrt{\rho} \sin(\varphi)$

Variants:

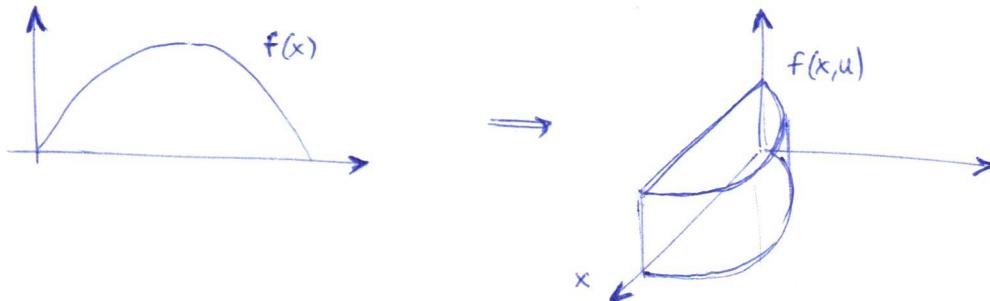
- $Y \sim N(\mu, \sigma^2)$: sample $X \sim N(0, 1) \Rightarrow Y = \sigma X + \mu$
- $\underline{Y} \sim N_k(\mu, \Sigma)$, $\Sigma = A^T A$: sample $\underline{X} \sim N_k(\underline{0}, I) \Rightarrow \underline{Y} = A \underline{X} + \mu$

FUNDAMENTAL THEOREM OF SIMULATIONS

Suppose we want to simulate from $X \sim f(x)$, where $f(x)$ is a density function defined on an arbitrary space.

Since $f(x) = \int_0^{f(x)} du \Rightarrow f(x)$ is the marginal of $(X, u) \sim U(\{(x, u) : u \leq f(x)\})$

It's a trick: augment the dimension to simplify the problem:



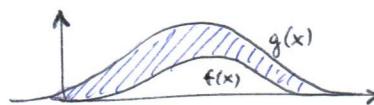
One method based on this theorem is for instance the Acceptance/Rejection method.

ACCEPTANCE / REJECTION METHOD

We are interested in sampling from $X \sim f(x)$.

We don't know how to do it so we use an instrumental random variable $Y \sim g(y)$.
 Y has to satisfy the following:

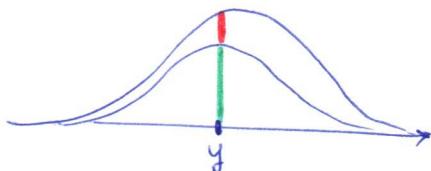
1. $S_x \leq S_y \leq 1R^k$ (support)
2. $f(x) \leq M g(x) \quad \forall x$



we sample using $g(\cdot)$ and we correct by the difference of the areas (///)

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 from $U | Y=y \sim U([0, Mg(y)])$ until $0 < U < f(y)$.

First we simulate on the support; obtaining y .



Then we simulate vertically (uniformly).
 If we fall in the acceptance zone we accept. If we fall in the rejection zone we reject

We augmented a dimension since we simulate twice (we sample horizontally and then vertically before accepting).

23/10/2020

3/5

```
# if not already installed
# install microbenchmark and RcppArmadillo
# about armadillo, take a look on:
# http://arma.sourceforge.net/docs.html

library(microbenchmark)
library(RcppArmadillo)
library(kernc)
library(ggplot2)
library(gridExtra)

#-----#
# # INVERSE TRANSFORM #
#-----#
```

```
#-----#
# Exponential((lambda))
#-----#
exp.custom <- function(n, lambda){
  results <- c()
  for(i in 1:n){
    u <- runif(1)
    results[i] <- -log(u) / lambda
  }
  return(results)
}
```

```
# plots comparison against the R base implementation
hist(rep(1000, 10), breaks = 50, freq = F, main = "", xlab = "")
lines(seq(0.001, 1, length.out = 100), dexp(seq(0.001, 1, length.out = 100), 10), col = 2)
```

```
#-----#
# c++ (using armadillo) implementation
cppFunction('arma::vec exp_custom_CPP(int n, double lambda){
  arma::vec results(n);
  double u;
  for(int i = 0; i < n; i++){
    u = arma::randu();
    results(i) = - log(u) / lambda;
  }
  return results;
}', depends = "RcppArmadillo")
```

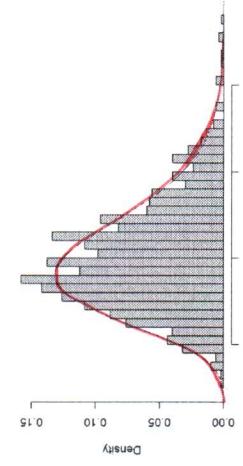
```
# comparison
microbenchmark(R = exp_custom(100, 1), CPP = exp_custom_CPP(100, 1), default = rexp(100, 1), times = 1000)
```

```
## Unit: microseconds
## expr   min   lq   mean median   uq   max neval
## R     198.5 207.65 303.0347 323.75 8228.9 1000
## CPP   7.0   7.98 12.97 8.65 10.10 1598.5 1000
## default 5.7   6.70 8.5299 7.20 8.05 96.4 1000
```

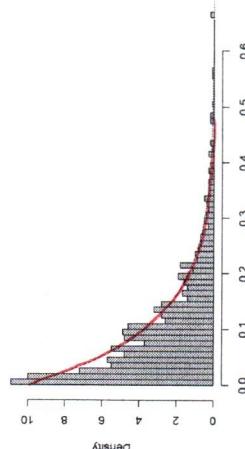
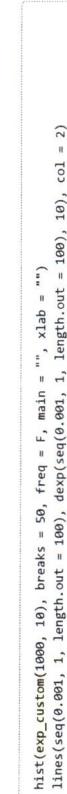
```
#-----#
# Gamma RV
#-----#
```

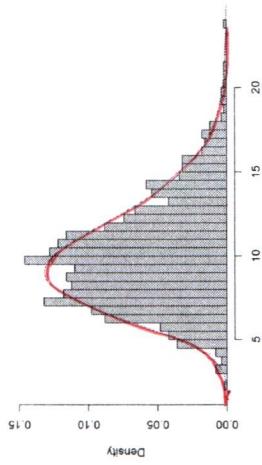
```
gamma.custom <- function(n, m, b){
  results <- c()
  for(i in 1:n){
    result.temp <- c()
    for(j in 1:m){
      u <- runif(1)
      result.temp[j] <- -log(u) / b
    }
    results[i] <- sum(result.temp)
  }
  return(results)
}
```

```
# plots
hist(gamma(1000, shape = 10, rate = 1), breaks = 50, freq = F, main = "", xlab = "")
lines(seq(0.001, 30, length.out = 1000), dgamma(seq(0.001, 30, length.out = 1000), shape = 10, rate = 1),
col = 2)
```



```
hist(gamma(custom(1000, 10, 1), breaks = 50, freq = F, main = "", xlab = ""))
lines(seq(0.001, 30, length.out = 1000), dgamma(seq(0.001, 30, length.out = 1000), shape = 10, rate = 1),
col = 2)
```





```

#-----#
# C++ (using armadillo) implementation
cppFunction("arma::vec Gamma_custom_CPP(int n, int m, double b){
    arma::double results(n);
    double u;
    arma::vec results_temp(m);
    for(int i = 0; i < n; i++){
        for(int j = 0; j < m; j++){
            u = arma::randu();
            results_temp(j) = - log(u) / b;
        }
        results(i) = arma::accu(results_temp);
    }
    return results;
}", depends = "RcppArmadillo")
}

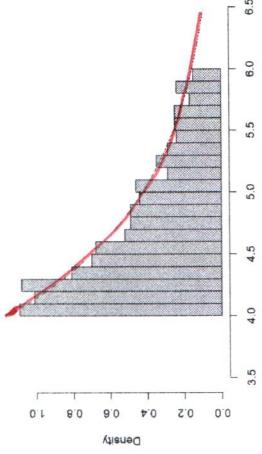
```

```

# -----#
# # Truncated RV
# -----#
# sample an exponential(1) on the subset of the support (8, 11)

trunc_exp.custom <- function(n, b, a1, a2){
  results <- c()
  for(i in 1:n) {
    u <- runif(1) * (exp(-b * a2) - exp(-b * a1)) + exp(-b * a1)
    results[i] <- rexp(1, 1/(log(u) / b + a1))
  }
}

```



```

# -----#
# c++ (using armadillo) implementation
cppFunction("arma::vec trunc_exp_custom_cpp(int n, double b, double a1, double a2){\n    arma::vec results(n);\n\n    double u;\n    for(int i = 0; i < n; i+){\n        u = arma::rand() * (exp(-b * a2) - exp(-b * a1)) + exp(-b * a1);\n        results(i) = -log(u) / b;\n    }\n\n    return results;\n}\n", depends = "RcppArmadillo")
# comparison\nmicrobenchmark(R = trunc_exp_custom(100, 1, 4, 6), CPP = trunc_exp_custom_CPP(100, 1, 4, 6), times = 10)

```

```

##           # Truncated RV
##           # sample an exponential(1) on the subset of the support (8, 11)

trunc_exp_custom <- function(n, b, a1, a2){
  results <- c()
  for(i in 1:n){
    u <- runif(1) * (exp(-b * a2) - exp(-b * a1)) + exp(-b * a1)
    results[i] <- -log(u) / b
  }
  return(results)
}

hist(trunc_exp_custom(1000, 1, 4, 6), freq = F, xlim = c(3.5, 6.5), breaks = 25, main = "", xlab = "")
curve(dexp(x, 1) / (rep(6, 1) - pexp(4, 1)), col = 2, add = T)

```

```

#-----#
# C++ (using armadillo) Implementation
cppFunction("arma::vec custom_sample_CPP(int n, arma::vec wei){
    arma::vec results(n);
    double u;
    arma::vec cumulate_wei = arma::cumsum(wei / accu(wei));
    for(int i = 0; i < n; i++){
        u = arma::randu();
        for(int j = 0; j < wei.n_elem; j++){
            if(u < cumulate_wei(j)){
                results(i) = j + 1;
                break;
            }
        }
    }
    return results;
}", depends = "RcppArmadillo")
}

microbenchmark(R = custom_sample_cpp(1000, c(0.1, 0.9, 0.1)),
               CPP = custom_sample_CPP(1000, c(0.1, 0.9, 0.1)),
               default = sample(1:3, size = 1000, replace = T, prob = c(0.1, 0.9, 0.1)), times = 1000)

## Unit: microseconds
## expr      min       lq      median       uq      max neval
## R 2494.7 2638.65 3392.8053 2829.00 2288.4 17432.6 1000
## CPP 22.4   24.78  49.2625  33.60  48.9   36966.7 1000
## default 22.6   25.60  38.1405  35.15  43.1   246.1  1000

```

```

#-----#
# SAMPLING A GAUSSIAN DISTRIBUTION
#-----#
norm_custom <- function(n, mean, sigma2){
  results <- matrix(0, ncol = 2, nrow = n/2)
  for(l in 1:(n/2)){
    u <- runif(2)
    rho = -2 * log(u[1])
    phi = 2 * pi * u[2]
    results[l,1] = sqrt(rho) * cos(phi)
    results[l,2] = sqrt(rho) * sin(phi)
  }
  return(as.vector(results) * sqrt(sigma2) + mean)
}

cppFunction('arma::vec norm_custom_CPP(int n, double mean, double sigma2){
  arma::vec results(n);
  double u1, u2, rho, phi;
  for(int i = 0; i < n/2; i++){
    u1 = arma::randu();
    u2 = arma::randu();
    rho = -2 * log(u1);
    phi = 2 * M_PI * u2;
    results[i,(n/2)] = sqrt(rho) * cos(phi);
    results[i,(n/2)] = sqrt(rho) * sin(phi);
  }
  return (results*sqrt(sigma2) + mean);
}', depends = "RcppArmadillo")

# plots

sample_R <- norm_custom_cpp(1000, 0, 1)
df1_R <- data.frame(y = sample_R)
df2_R <- data.frame(x = seq(-4, 4, length.out = 100),
                     dens = dnorm(seq(-4, 4, length.out = 100), 0, 1))
df3_R <- data.frame(cmean = cumsum(sample_R) / 1:length(sample_R),
                     x = 1:length(sample_R))

p1 <- ggplot() +
  geom_histogram(data = df1_R, mapping = aes(y, stat(density)), alpha = 0.3, col = 1, bins = 30) +
  geom_line(data = df2_R, mapping = aes(x, y = dens), col = 2+
  theme_bw()) +
  ggtitle("R implementation")

l1 <- ggplot() +
  geom_line(data = df3_R, mapping = aes(x = x, y = cmean)) +
  geom_line(data = df3_R, mapping = aes(x = x, y = mean(sample_R)), col = 2+
  theme_bw())

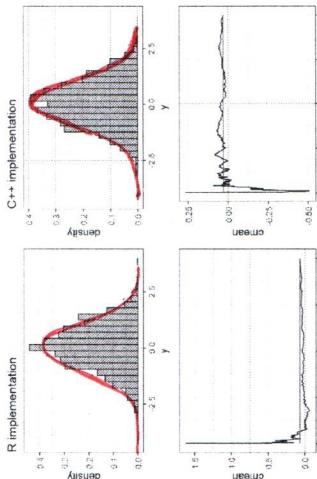
sample_CPP <- norm_custom_CPP(1000, 0, 1)
df1_CPP <- data.frame(y = sample_CPP)
df2_CPP <- data.frame(x = seq(-4, 4, length.out = 100),
                      dens = dnorm(seq(-4, 4, length.out = 100), 0, 1))
df3_CPP <- data.frame(cmean = cumsum(sample_CPP) / 1:length(sample_CPP),
                       x = 1:length(sample_CPP))

p2 <- ggplot() +
  geom_histogram(data = df1_CPP, mapping = aes(y, stat(density)), alpha = 0.3, col = 1, bins = 30) +
  geom_line(data = df2_CPP, mapping = aes(x, y = dens), col = 2+
  theme_bw()) +
  ggtitle("C++ implementation")

l2 <- ggplot() +
  geom_line(data = df3_CPP, mapping = aes(x = x, y = cmean)) +
  geom_line(data = df3_CPP, mapping = aes(x = x, y = mean(sample_CPP)), col = 2 +
  theme_bw())

grid.arrange(p1, p2, l1, l2, nrow = 2)

```

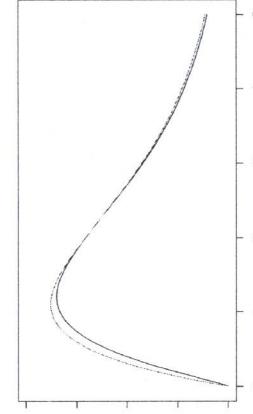


卷之三

```
# Comparison
microbenchmark(R = norm_custom(100, 0, 1),
               CPP = norm_custom_CPP(100, 0, 1),
               default = rnorm(100, mean = 0, sd = 1),
               times = 100)
```

unit: microseconds		expr			min	1q	mean	median	uq	max	neval
R	124.0	129.5	145.7	163.15	22024.4	100					
CPP	10.5	11.25	15.292	11.9	13.25	150	1.1				
default	9.8	10.35	12.927	10.9	11.90	91.1	1.1				

```
    d(digamma(x, shape = 2.2, rate = 1), xlab = "x")
    d(c(1, 2) * digamma(x, shape = 2, rate = 2/2.2), add = T, col = 2)
```



```

curve(dgamma(x, shape = 2.2, rate = 1), xlim = c(4, 5), xlab = "", ylab = "")
curve(c * dgamma(x, shape = 2, rate = 2/2.2), add = T, col = 2)

```

```
a = round(a) * exp(-(a - r0))  
b = "", ylab = "")
```

```
#-----#
gamma_AR <- function(n, a, b) {
  acc <- 0
  k <- 0
  a_star <- round(a)
```

```

M <- gamma(a_star) / gamma(a) * b_star^(a - a_star) / ((a - a_star)^(a - b_star)) * exp(-
(a - a_star))

results <- c()
while(acc < c()){
  y <- rgamma(1, shape = a_star, rate = b_star)
  u <- runif(1)
  if(u < dgamma(y, shape = a, rate = b) / (M * dgamma(y, shape = a_star, rate = b_star))){
    acc = acc + 1
    results <- c(results, y)
  }
  k <- k + 1
}

print(paste0("Acc. ratio = ", round(acc / k, digits = 3)))
return(results)
}

```

卷之三

卷之三

```

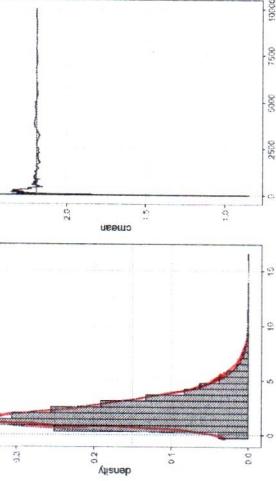
df1 <- data.frame(y = sample_GAMMA)
df2 <- data.frame(x = seq(0, 15, length.out = 100),
                  dens = dgamma(seq(0, 15, length.out = 100), shape = 2.2, rate = 1))
df3 <- data.frame(cumsum.sample_GAMMA = cumsum(sample_GAMMA) / 1:length(sample_GAMMA),
                  x = 1:length(sample_GAMMA))

p <- ggplot() +
  geom_histogram(data = df1, mapping = aes(y, stat(density)), alpha = 0.3, col = 1, bins = 30) +
  geom_line(data = df2, mapping = aes(x = x, y = dens), col = 2) +
  theme_bw() +
  ggtitle("GAMMA(2.2, 1)")

l <- ggplot() +
  geom_histogram(data = df3, mapping = aes(x = x, y = cmean)) +
  geom_line(data = df3, mapping = aes(x = x, y = mean(sample_GAMMA)), col = 2) +
  theme_bw() +
  ggtitle("")

grid.arrange(p, l, nrow = 1)

```



```
#> -----
#> RcppFunction("arma::vec gamma_AR_CPP(int n, double a, double b){
#>   arma::vec results(n, arma::fill::zeros);
#>   int k = 0;
#>   double y, u;
```

```
double a_star = round(a);
double b_star = a_star / a;
double M = pow(b_star, -a_star) * pow((a - a_star) / (1 - b_star), (a - a_star)) *
exp(-(a - a_star) + igamma(a_star) - igamma(a));
while(acc < n - 1){
y = arma::randg(a_star);
u = arma::rand();
if(u < pow(b_star, -a_star) * distri_param(a, 1 / b));
exp(-y * (1 - b_star) + igamma(a_star) - igamma(a)) / M);
results(acc) = y;
acc += 1;
}
k += 1;
}
```

```
Rcpp::Rcout << (double) acc / (double) k;
return (results);
}","depends = "RcppArmadillo")
```

```
microbenchmark(R = Gamma.AR(100, 2, 2, 1),
CPP = Gamma.AR_CPP(100, 2, 2, 1),
times = 1000)
```

```
## [1] "Acc. ratio = 0.952"
## [1] "Acc. ratio = 0.952"
## 0.933962.9519230.970588[1] "Acc. ratio = 0.981"
## 0.961659.970588[1] "Acc. ratio = 0.926"
## .
```

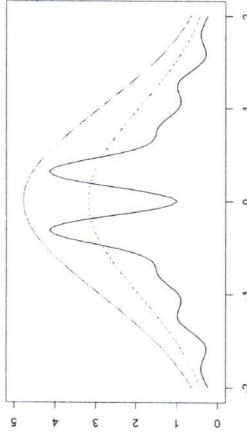
```
## Unit: microseconds
## expr min    q50   median    q75   max  nrow
## exp(x) 10 101.85 151.63743 118.85 1393.25 10000
## CPP 339 363.35 487.6909 385.95 415.66 1655.6 10000
```

```
#> -----
#> # ACCEPTANCE-REJECTION ~~~~ 2
```

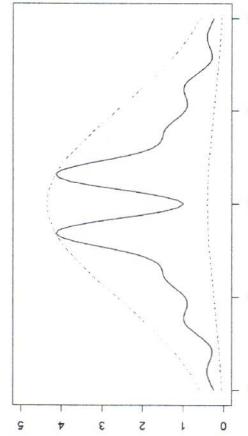
```
fx <- function(x) exp(-x^2 / 2) * (sin(6 * x)^2 + 3 * cos(x)^2 * sin(4*x)^2 + 1)
integrate(fx, -Inf, Inf)$value
```

```
## [1] 5.89434
```

```
curve(fx, xlim = c(-2, 2), ylim = c(0, 5), xlab = "", ylab = "")
curve(dnorm(x)*8, col = 2, lty = 2, add = 1)
curve(dnorm(x)*12, col = 2, lty = 1, add = 1)
```



```
M <- optim(theta, fn = function(x) -fx(x) / dnorm(x), method = "Brent",
            curve(fx, xlim = c(-2, 2), ylim = c(0, 5), xlab = "", ylab = ""))
curve(dnorm(x), col = 2, lty = 2, add = 1)
curve(dnorm(x)*M, col = 3, lty = 2, add = 1)
```



```
fx_AR <- function(n, M){
  acc < 0
  k < 0
  results <- c()
  while(acc < n){
    y <- rnorm(1)
    u <- runif(1)
    if(u < fx(y) / (M * dnorm(y))){
      acc = acc + 1
      results <- c(results, y)
    }
  }
  k <- k + 1
}
print(paste0("Acc. ratio = ", round(acc / k, digits = 3)))
return(results)
}

# plots
fx_sample <- fx_AR(10000, M)
```

```
## [1] "Acc. ratio = 0.546"
```

```

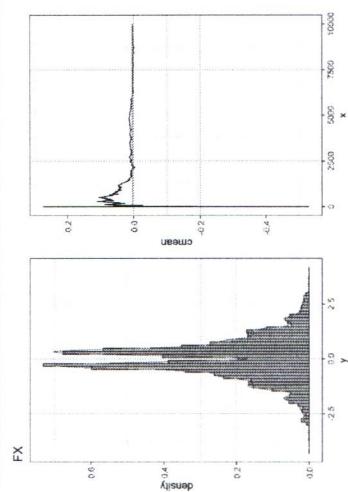
df1 <- data.frame(y = fx_sample)
df2 <- data.frame(x = seq(-4, 4, length.out = 1000),
  dens = fx$seq(-4, 4, length.out = 1000) / integrate(fx, -Inf, Inf)$value)
df3 <- data.frame(mean = sapply(1:length(fx_sample), function(x) mean(fx_sample[1:x])),
  x = 1:length(fx_sample))

p <- ggplot() +
  geom_histogram(data = df1, mapping = aes(y, stat(density)), alpha = 0.3, col = 1, bins = 60) +
  geom_line(data = df2, mapping = aes(x = x, y = dens), col = 2) +
  theme_bw() +
  ggtitle("fx")

l <- ggplot() +
  geom_line(data = df3, mapping = aes(x = x, y = mean)) +
  geom_line(data = df3, mapping = aes(x = x, y = mean(fx_sample)), col = 2) +
  theme_bw() +
  ggtitle("mean")

grid.arrange(p, l, nrow = 1)

```



ACCEPTANCE/REJECTION METHOD

29/10/2020

1/4

We want to sample from $X \sim f_x$.

Prop. Let $X \sim f_x$ and $Y \sim g_x$. s.t.:

$$1. Sf_x \subseteq Sg_x \subseteq \mathbb{R}^K$$

$$2. \exists M > 0 : f_x \leq Mg_x$$

→ To simulate from X is sufficient to simulate from $Y \sim g_x$ and $U | Y=y \sim U([0, Mg_x(y)])$ until $0 < U < f(y)$

The proposition is saying: $\Pr(Y \leq x | U \leq \frac{f_x(Y)}{Mg_x(Y)}) = \Pr(X \leq x)$

proof. ($u = u$)

$$\begin{aligned} \Pr(Y \leq x | U \leq \frac{f_x(Y)}{Mg_x(Y)}) &= \frac{\Pr(Y \leq x, U \leq \frac{f_x(Y)}{g_x(Y) \cdot M})}{\Pr(U \leq \frac{f_x(Y)}{Mg_x(Y)})} \\ &= \frac{\mathbb{E}_Y [\Pr(Y \leq x, U \leq \frac{f_x(Y)}{Mg_x(Y)} | Y=y)]}{\mathbb{E}_Y [\Pr(U \leq \frac{f_x(Y)}{Mg_x(Y)} | Y=y)]} \\ &= \frac{\int \Pr(Y \leq x, U \leq \frac{f_x(Y)}{Mg_x(Y)} | Y=y) \cdot g(y) dy}{\int \Pr(U \leq \frac{f_x(Y)}{Mg_x(Y)} | Y=y) g(y) dy} \\ &= \frac{\int \mathbb{1}_{(-\infty, x)}(y) \Pr(U \leq \frac{f_x(y)}{Mg_x(y)}) g(y) dy}{\int \Pr(U \leq \frac{f_x(y)}{Mg_x(y)}) g(y) dy} \\ &\stackrel{\text{Here } U \sim U([0,1]) \text{ (because of the condition), we use its cdf}}{=} \frac{\int \mathbb{1}_{(-\infty, x)}(y) \frac{f_x(y)}{Mg_x(y)} g(x) dy}{\int \frac{f_x(y)}{Mg_x(y)} g(x) dy} \\ &= \frac{\int \mathbb{1}_{(-\infty, x)}(y) F_x(y) dy}{\int F_x(y) dy} = 1 \end{aligned}$$

Algorithm:

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

1. Sample $Y \sim g_x$

2. Sample $U \sim U([0,1])$

3. If $U \leq \frac{f(y)}{Mg(y)}$ set $X_i = Y$, otherwise go back to 1.

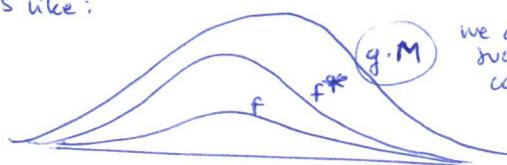
Remark 1. The algorithm works also with not normalized densities.

Suppose to have: $f_x = \alpha f_x^*$, $g_x = \beta g_x^*$, where f_x^* and g_x^* are not normalized densities.

$$\Rightarrow f_x = Mg_x \iff f_x^* = M^* g_x^*$$

we can always choose another M^* (s.t. the algorithm works)

It's like:



we can always find an M such that $M \cdot g$ (or in case $M^* \cdot g^*$) covers both f and f^*

Remark 2. What is the probability of accepting a proposed value?

$$\begin{aligned} \Pr(U \leq \frac{f_x(y)}{Mg_x(y)}) &= \int \Pr(U \leq \frac{f_x(y)}{Mg_x(y)} | Y=y) g(y) dy \\ &= \int \frac{f_x(y)}{Mg_x(y)} g(y) dy \\ &= \frac{1}{M} \int f_x(y) dy = \frac{1}{M} \end{aligned}$$

if we want to optimize the algorithm
we have to maximize this probability
(and if we have to reduce M)



Example.

We want to sample from $f_x \stackrel{d}{=} \text{Gamma}(\alpha, 1)$, $\alpha > 1$.

We use as proposal $g_x \stackrel{d}{=} \text{Gamma}(m, b)$, with $m = \lfloor \alpha \rfloor$.

Find M (if \exists) s.t. $f_x \leq Mg_x$.

We study:

$$\frac{f_x}{g_x} = \frac{\frac{1}{\Gamma(\alpha)} x^{\alpha-1} e^{-x}}{\frac{1}{\Gamma(m)} b^m x^{m-1} e^{-bx}} = \frac{\Gamma(m)}{\Gamma(\alpha)} b^{-m} \underbrace{x^{\alpha-m} e^{-(1-b)x}}_{\text{Gamma}(\alpha-m+1, 1-b)}$$

(Conditions)

$$\text{If } 1-b > 0 \Rightarrow b < 1$$

$$\text{If } \alpha - m + 1 > 0 \Rightarrow \alpha - \lfloor \alpha \rfloor - 1 > 0 \quad (\text{always})$$

The only condition is $b \in (0, 1)$ (to guarantee a proper density).

For $b \in (0, 1)$ the maximum of the Gamma is in its mode:

$$\text{mode}(\text{Gamma}(\alpha-m+1, 1-b)) = \frac{\alpha - \lfloor \alpha \rfloor}{1-b} = x \text{ in } *$$

\Rightarrow the maximum value of the ratio is in $\frac{\alpha - \lfloor \alpha \rfloor}{1-b}$ (ratio: $\frac{f(\cdot)}{g(\cdot)}$)

$$\frac{f_x}{g_x} \leq \frac{\Gamma(\lfloor \alpha \rfloor)}{\Gamma(\alpha)} b^{-\lfloor \alpha \rfloor} \left(\frac{\alpha - \lfloor \alpha \rfloor}{1-b} \right)^{\alpha - \lfloor \alpha \rfloor} e^{-(1-b)\left(\frac{\alpha - \lfloor \alpha \rfloor}{1-b}\right)} := M(b)$$

since α is fixed, the expression is a function of b

We can optimize it in terms of b .

Notice: it's simpler to study $M^{-1}(b)$: (not the inverse, $\frac{1}{M(b)}$)

We have that: $\frac{1}{M(b)} \propto b^{\lfloor \alpha \rfloor} (1-b)^{\alpha - \lfloor \alpha \rfloor} \propto \text{Beta}(\lfloor \alpha \rfloor + 1, \lfloor \alpha \rfloor - \alpha + 1)$

Optimize in terms of b = minimize $M(b)$ = maximize $M^{-1}(b)$.

To maximize $\frac{1}{M(b)} \stackrel{d}{=} \text{Beta}(\cdot, \cdot) \rightarrow$ we consider the mode of $\text{Beta}(\cdot, \cdot)$

The mode of $\text{Beta}(\cdot, \cdot)$ \exists only if both powers are > 1 .

The mode of $\text{Beta}(\cdot, \cdot)$ is in $\frac{\lfloor \alpha \rfloor}{2}$

$$\Rightarrow b_{\text{optimal}} = \frac{\lfloor \alpha \rfloor}{\alpha} \Rightarrow M_{\text{optimal}} = M(b_{\text{optimal}})$$

IMPORTANCE SAMPLING

Let f_x be a (regular enough) density and h_x a measurable function.

We're interested in: $E_f[h(X)]$?

$$E_f[h(X)] = \int_{\mathbb{R}} h(x) f(dx)$$

A first trivial approach: Monte Carlo method:

$$E_f[h(X)] \approx \frac{1}{m} \sum_{i=1}^m h(x_i) \quad \text{with } x_1, \dots, x_m \stackrel{iid}{\sim} f_x$$

By the (strong) law of large numbers ((S)LLN):

$$\frac{1}{m} \sum_{i=1}^m h(x_i) \xrightarrow{\text{a.s.}} E_f[h(X)]$$

However, in some scenarios we're not able to sample from f_x . To overcome the problem we introduce g_x , an instrumental random variable.

- We generate $x_1, \dots, x_m \sim g_x$
- $E_f[h(X)] \approx \frac{1}{m} \sum_{i=1}^m \frac{f(x_i)}{g(x_i)} h(x_i)$

We introduce the weight function: $w(x) = \frac{f(x)}{g(x)}$ ($w(x_i) = \frac{f(x_i)}{g(x_i)}$)

We have that:

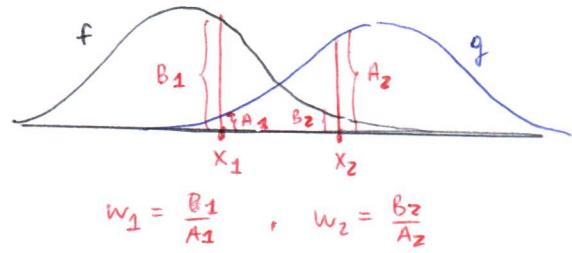
$$\frac{1}{m} \sum_{i=1}^m \frac{f(x_i)}{g(x_i)} h(x_i) \xrightarrow{\text{a.s.}} E_g\left[\frac{f(x)}{g(x)} h(x)\right] = \int \frac{f(x)}{g(x)} h(x) g(x) dx \\ = \int f(x) h(x) dx = E_f[h(X)]$$

Remark: we can use importance sampling to sample from a distribution.

We can approximate $f(x)$:

$$\begin{aligned} \tilde{f}(x) &= \frac{1}{m} \sum_{i=1}^m \frac{f(x_i)}{g(x_i)} \delta_{x_i}(x) \\ &\equiv \frac{1}{m} \sum_{i=1}^m w(x_i) \delta_{x_i}(x) \end{aligned}$$

→ when we're not able to sample from $f(x)$ we can sample from $g(x)$ and re-weight the probability by $w(x)$.



Algorithm:

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

1. Sample $x_1, \dots, x_m \sim g$ (m large → good (global) sample of f)
2. Compute (w_1, \dots, w_m) with $w_j = \frac{f(x_j)}{g(x_j)}$
3. Sample from x_1, \dots, x_m with probabilities (w_1, \dots, w_m)

Theorem. The optimal choice of g_x (that minimizes the variance of the importance sampling estimator) is:

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

proof.

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

$$\text{Notice that: } E_g\left[h(x) \frac{f(x)}{g(x)}\right] = \int h(x) \frac{f(x)}{g(x)} g(x) dx = \int h(x) f(x) dx \quad \text{if } g(x) = f(x)$$

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

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

$$f(\cdot), g(\cdot) \geq 0$$

with this we minimize $E_g[(\dots)^2]$

Suppose we want to evaluate the tails of a distribution:

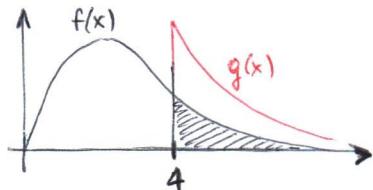
Example.

$$X \sim \text{Weibull}(2, 2)$$

We want to evaluate: $P(X > 4)$?

$$f(x) = \frac{k}{\lambda} x^{k-1} e^{-\left(\frac{x}{\lambda}\right)^k}$$

We want to evaluate:



We want to treat it as an importance sampling problem:

$$f(x) \sim \text{Weibull}(2, 2)$$

$$h(x) = \mathbb{1}_{(4, \infty)}(x)$$

$$\underline{g(x) \sim \mathcal{E}(\lambda) + 4}$$

```

library(microbenchmark)
library(Rcpp)
library(RcppArmadillo)
library(Rcpp)
library(ggplot2)
library(gridExtra)
library(ggpubr)

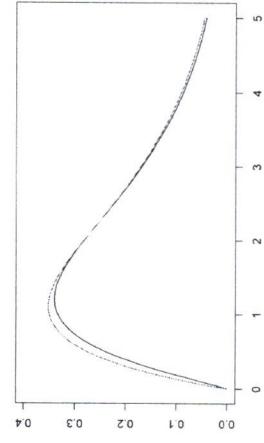
#--- ACCEPTANCE / REJECTION : gamma ( . )
#---



a <- 2.2
b <- round(a) / a
C = gamma(round(a)) / gamma(a) * b^(round(a)) * ((a - round(a))^(1 - b))^(a - round(a)) / (1 - b)^(a - round(a))
und(a))

curve(dgamma(x, shape = 2.2, rate = 1), xlim = c(0, 5), ylim = c(0, 0.4), xlab = "", ylab = "")
curve(C * dgamma(x, shape = 2, rate = 2), add = T, col = 2)

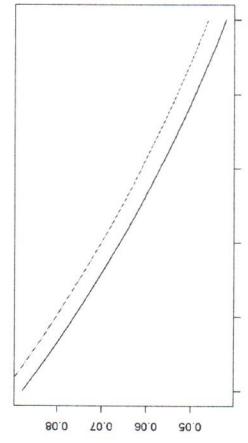
```



```

curve(dgamma(x, shape = 2.2, rate = 1), xlim = c(4, 5), xlab = "", ylab = "")
curve(C * dgamma(x, shape = 2, rate = 2/2.2), add = T, col = 2)

```



```

## [1] "Acc. ratio = 0.951"

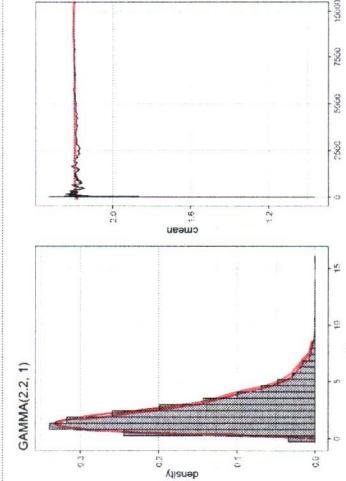
df1 <- data.frame(y = sample_GAMMA)
df2 <- data.frame(x = seq(0, 15, length.out = 100),
                  dens = dgamma(seq(0, 15, length.out = 100), shape = 2.2, rate = 1))
df3 <- data.frame(cmean = supply(1:length(sample_GAMMA), function(x) mean(sample_GAMMA[1:x])), x = 1:length(sample_GAMMA))

p <- ggplot() +
  geom_histogram(data = df1, mapping = aes(y, stat(density)), alpha = 0.3, col = 1, bins = 30) +
  geom_line(data = df2, mapping = aes(x = x, y = dens), col = 2) +
  theme_bw() +
  ggtitle("GAMMA(2, 2, 1)")

l <- ggplot() +
  geom_line(data = df3, mapping = aes(x = x, y = cmean)) +
  geom_line(data = df3, mapping = aes(x = x, y = mean(sample_GAMMA)), col = 2) +
  theme_bw() +
  ggtitle("")

grid.arrange(p, l, nrow = 1)

```



```

#-----#
cppFunction("arma::vec gamma_AR_CPP(int n, double a, double b){
  arma::vec results(n, arma::fill::zeros);
  int acc = 0;
  int k = 0;
  double y, u;
  double a_star = round(a);
  double b_star = a_star / a;
  double M = pow(b_star, -a_star) * pow((a - a_star) / (1 - b_star)), (a - a_star)) *
    exp(-(a - a_star) + lgamma(a_star));
  while(acc < n - 1){
    u = arma::randu();
    y = arma::randn();
    if(u < pow(b_star, -a_star) * pow(y, a - a_star) *
      exp(-y * (1 - b_star) + lgamma(a_star) - lgamma(a)) / M){
      results(acc) = y;
      acc += 1;
    }
  }
  Rcpp::Rcout << (double) acc / (double) k;
  return (results);
} ", depends = "RcppArmadillo")
microbenchmark(gamma_AR_CPP(100, 2, 1),
  gamma_AR_CPP(100, 2, 1),
  times = 100)
}

## [1] "Acc. ratio = 0.935"
## # 0.96165[1] "Acc. ratio = 0.962"
## # 0.9801988. 951923[1] "Acc. ratio = 0.935"
## ..

## Unit: microseconds
## expr   min   q1   median   q3   max  neval
## gamma_AR(100, 2, 1) 1938.961 1190.958 1676.8950 1522.052 1712.3510
## gamma_AR_CPP(100, 2, 1) 344.101 374.7015 434.1711 408.751 431.9695
## ..
```

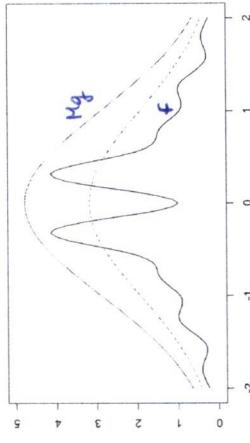
```

## # Unit: microseconds
## # ACCEPTANCE-REJECTION ---- 2
## #-----#
fx <- function(x) exp(-x^2 / 2) * (sin(6 * x)^2 + 3 * cos(x)^2 * sin(4*x)^2 + 1)
integrate(fx, -Inf, Inf)$value
y <- rnorm(1)
u <- runif(1)
if(u <= fx(y) / (M * dnorm(y))){
  acc = acc + 1
  results <- c(results, y)
}
k <- k + 1
print(paste0("Acc. ratio = ", round(acc / k, digits = 3)))
return(results)
}
```

```

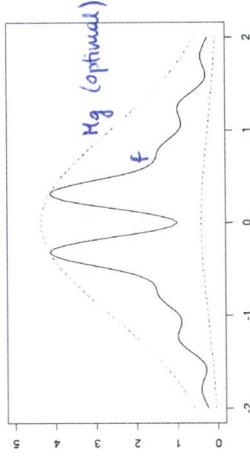
## [1] 5.89434
curve(fx, xlim = c(-2, 2), col = 2, lty = 2, add = T)
curve(dnorm(x)*8, col = 2, lty = 2, add = T)
curve(dnorm(x)*12, col = 2, lty = 1, add = T)

## [1] "Acc. ratio = 0.535"
```



```

M <- optim(theta, fn = function(x) -fx(x) / dnorm(x), method = "Brent",
  lower = -10, upper = 10)$value
```



```

fx_AR <- function(n, M){
  acc <- 0
  k <- 0
  results <- c()
  while(acc < n){
    y <- rnorm(1)
    u <- runif(1)
    if(u <= fx(y) / (M * dnorm(y))){
      acc = acc + 1
      results <- c(results, y)
    }
    k <- k + 1
  }
  print(paste0("Acc. ratio = ", round(acc / k, digits = 3)))
  return(results)
}

fx_sample <- fx_AR(10000, M)

## [1] "Acc. ratio = 0.535"
```

29/10/2020

4/4

```

#-----#
# dummy sampler: counts how many values are greater than a threshold
dummy_wei <- function(n, a, b, c){
  temp <- rweibull(n, a, b)
  return(mean(temp > c))
}

#-----#
# importance sampling: sample from an exponential distribution translated by the threshold
imp_wei <- function(n, a, b, c){
  we1 <- c()
  for(i in 1:n){
    temp <- rgamma(1, 1, 1) + c
    we1[i] <- dweibull(temp, a, b) / dgamma(temp - c, 1, 1)
  }
  return(mean(we1))
}

#-----#
# replicate the evaluation of the tail one hundred of times
results <- rep(0, 200)
for(i in 1:100){
  results[i] <- dummy_wei(100, 2, 2, 4)
  results[i+100] <- imp_wei(100, 2, 2, 4)
}

df <- data.frame(res = results, type = rep(c("DUMMY", "IS-EXP"), each = 100))

p1 <- ggplot(df) +
  geom_boxplot(aes(x = type, y = res, fill = type), notch=TRUE) +
  geom_hline(yintercept = pweibull(4, 2, 2, lower.tail = FALSE), col = 1, lwd = 1, lty = 2) +
  theme_bw() +
  ylab("value") +
  xlab("type")

p2 <- ggplot(df) +
  geom_boxplot(aes(x = type, y = res, fill = type)) +
  geom_hline(yintercept = pweibull(4, 2, 2, lower.tail = TRUE), col = 1, lwd = 1, lty = 2) +
  theme_bw() +
  ylab("value") +
  xlab("type")

p3 <- ggplot(df) +
  geom_boxplot(aes(x = type, y = res, fill = type)) +
  geom_hline(yintercept = pweibull(4, 1000, 2, 2, 4)) +
  theme_bw() +
  ylab("value") +
  xlab("type")

p4 <- ggplot(df) +
  geom_boxplot(aes(x = type, y = res, fill = type)) +
  geom_hline(yintercept = pweibull(4, 2, 2, lower.tail = TRUE), col = 1, lwd = 1, lty = 2) +
  theme_bw() +
  ylab("value") +
  xlab("type")

```

```

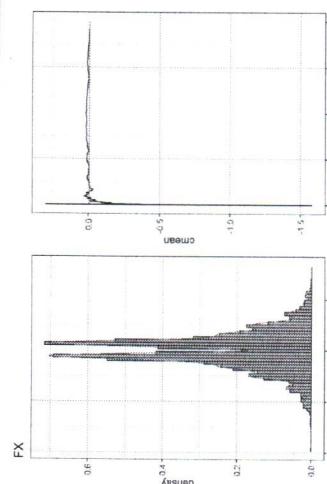
df1 <- data.frame(y = fx_sample)
df2 <- data.frame(x = seq(-4, 4, length.out = 1000),
  dens = fx$seq[-4:4], length.out = 1000) / integrate(fx, -inf, Inf)$value,
  x = 1:length(fx_sample))

p <- ggplot() +
  geom_histogram(data = df1, mapping = aes(y, stat(density)), alpha = 0.3, col = 1, bins = 60) +
  geom_line(data = df2, mapping = aes(x = x, y = dens), col = 2) +
  theme_bw() +
  ggtitle("FX")

l <- ggplot() +
  geom_line(data = df3, mapping = aes(x = x, y = mean)) +
  geom_line(data = df3, mapping = aes(x = x, y = mean(fx.sample)), col = 2) +
  theme_bw() +
  ggtitle("")

grid.arrange(p, l, nrow = 1)

```



Right tail of a Weibull(2,2)

evaluate the probability on the right tail of a weibull(2,2) distribution

```

curve(dweibull(x, 2, 2), xlim = c(0, 8), ylim = c(0, 1), main = "", ylab = "dens")
lines(seq(4, 10, length.out = 100), dexp(seq(0, 10, length.out = 100)), col = 4)
lines(seq(4, 10, length.out = 100), dexp(seq(4, 10, length.out = 100) - 4), col = 2)

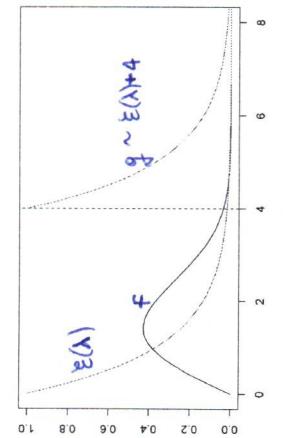
```

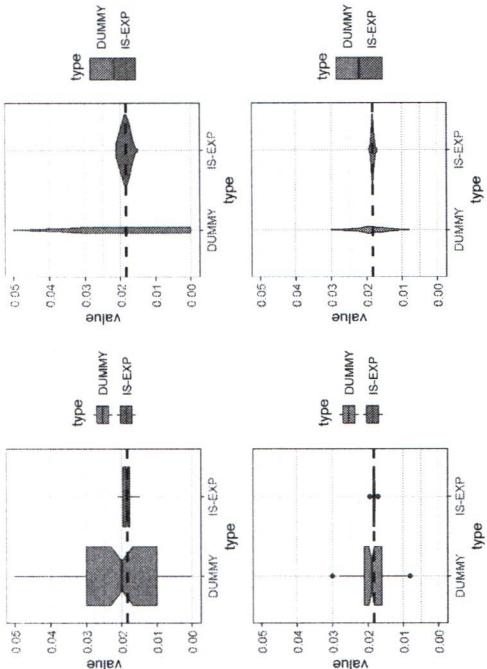
#

IMPORTANCE SAMPLING

#

curve(dweibull(x, 2, 2), xlim = c(0, 8), ylim = c(0, 1), main = "", ylab = "dens")
lines(v = 4, lty = 2)





MCMC

We want to sample from $f_X(x)$ with support S_X .

We use a sequence $\{x^{(t)}\}$ with ergodic distribution equivalent to $f_X(x)$.

How to generate this sequence $(\{x^{(t)}\})$?

A possible strategy is the use of the METROPOLIS-HASTINGS ALGORITHM.

We use a proposal distribution to propose a candidate at the t -step:

$$q(y|x^{t-1}) = q(y|x) \quad (\text{notations}).$$

We need some requirements:

- $q(\cdot|x)$ is a density $\forall x \in S_X$

- $q(y|\cdot)$ is measurable $\forall y$

- $S_X \subseteq S_q \subseteq \mathbb{R}^K$ (the support of the proposal must cover the support of the target var)

How does it work?

We propose a value for $f(x)$ by the conditional distribution $q(y|x)$ and we have a probability of accepting/rejecting:

$$\alpha(y,x) = \min \left\{ 1, \frac{f(y) q(x|y)}{f(x) q(y|x)} \right\}$$

$\Rightarrow \begin{cases} 1. \text{ We propose a value } y \sim q(y|x) \\ 2. \text{ We accept with probability } \alpha(x,y) \end{cases}$ (= probability to move from x to y , while $1 - \alpha(y,x)$ is the probability of staying in x)

Algorithm:

- Set $t=0$, $y^0 = x_0$ = initial value
- for $t=1, \dots, n_{\text{rep}}$:
 1. Sample $Y \sim q(y|x^{t-1})$
 2. Sample $U \sim U([0,1])$
 3. If $U \leq \alpha(x,y)$ $\Rightarrow x^t = Y$
otherwise $\Rightarrow x^t = x^{t-1}$

Theorem.

Let $\{x^{(t)}\}$ be a chain produced by a Metropolis-Hastings algorithm.

1. The kernel associated with the chain satisfies the "detailed balanced condition":

$$\exists f: k(y,x) f(y) = k(x,y) f(x)$$

2. f is the ergodic distribution of the chain

About the kernel:

$$k(y,x) = \alpha(y,x) q(y|x) + r(x) \delta(x)$$

$$r(x) = 1 - \int \alpha(y,x) q(y|x) dy$$

Remark 1. We can use this strategy to sample from the posterior distribution

$$f(\theta) = \pi(\theta | x_1, \dots, x_n)$$

Remark 2. It works also up to a constant:

$$\frac{f(y) q(x|y)}{f(x) q(y|x)} = \frac{M f^*(y) q(x|y)}{M f^*(x) q(y|x)}$$

f^* = not normalized function

= we can sample from not normalized functions
(In fact recall that, instead of computing exactly the posterior distribution we just have:

$$\pi(\theta | x) \propto L(\theta, x) \pi(\theta) \quad \leftarrow \text{not normalized}$$

We can tune the algorithm choosing different $q(y|x)$.

For instance:

- Independent Metropolis-Hastings : $q(y|x) \stackrel{d}{=} q(y) \quad (\perp x)$
- Random-Walk Metropolis-Hastings : $q(y|x) \stackrel{d}{=} q(y-x)$

Example 1.

Suppose to have $Y_1, \dots, Y_n | \mu, \lambda \stackrel{iid}{\sim} \text{Cauchy}(\mu, \lambda)$:

$$f_Y(y) = \frac{e^{-\lambda}}{\pi(1 + (y - \mu)^2 e^{-2\lambda})} \stackrel{\perp}{\sim} (-\infty, \infty) \quad \lambda, \mu \in \mathbb{R}$$

We want to make posterior inference on μ and λ .

We assume as prior:

$$\pi(\mu, \lambda) = N(\mu, m_1 \sigma_1^2) \cdot N(\lambda, m_2 \sigma_2^2)$$

The Cauchy is "the bad guy" because: $\nexists E[Y]$, $\nexists \text{Var}(Y)$.

Posterior distribution:

$$\begin{aligned} \pi(\mu, \lambda | Y_1, \dots, Y_n) &\propto \prod_{i=1}^n f_Y(y_i) \pi(\lambda) \pi(\lambda) \\ &\propto \frac{e^{-n\lambda}}{\pi^n \cdot \prod_{i=1}^n (1 + (y_i - \mu)^2 e^{-2\lambda})} \pi(\mu) \pi(\lambda) \end{aligned}$$

We define a MH algorithm using as proposal:

$$q(\mu^*, \lambda^* | \mu^{t-1}, \lambda^{t-1}) = N_2\left(\begin{pmatrix} \mu^{t-1} \\ \lambda^{t-1} \end{pmatrix}, \Sigma\right)$$

$$\alpha(\mu^*, \lambda^*, \mu^{t-1}, \lambda^{t-1}) = \frac{\pi(\mu^*, \lambda^* | Y)}{\pi(\mu^{t-1}, \lambda^{t-1} | Y)} \frac{q(\mu^{t-1}, \lambda^{t-1} | \mu^*, \lambda^*)}{q(\mu^*, \lambda^* | \mu^{t-1}, \lambda^{t-1})}$$

$$\begin{aligned} \log(\alpha(-, -)) &= \log(\pi(\mu^*, \lambda^* | Y)) - \log(\pi(\mu^{t-1}, \lambda^{t-1} | Y)) \\ &\stackrel{!}{=} \sum_{i=1}^n \log f(y_i | \mu^*, \lambda^*) + \log \pi(\mu^*) + \log \pi(\lambda^*) + \\ &\quad - \sum_{i=1}^n \log f(y_i | \mu^{t-1}, \lambda^{t-1}) - \log \pi(\mu^{t-1}) - \log \pi(\lambda^{t-1}) \end{aligned}$$

In case of random walk Metropolis-Hastings we can cancel this because we have a symmetric distribution

Example 2.

Sample of a mixture of 2 gaussian distributions.

$$f(x | \mu_1, \mu_2, \sigma_1^2, \sigma_2^2, \beta) \stackrel{d}{=} \beta N(\mu_1, \sigma_1^2) + (1-\beta) N(\mu_2, \sigma_2^2)$$

We want to use as proposal:

$$q(y | x^{t-1}) \stackrel{d}{=} N(y; x^{t-1}, \sigma_q^2)$$

We'll consider: $\beta = 0.2$, $\mu_1 = -2.5$, $\mu_2 = 2.5$, $\sigma_1^2 = \sigma_2^2 = 1$

ADAPTIVE METROPOLIS HASTINGS

Instead of considering a proposal: $q(y|x)$
 we consider a family of proposals: $\{q_\gamma(y|x)\}_{\gamma \in \Gamma}$.

In this way the proposal may change
 during the simulation

Associated with the family of proposals we have a family of transition kernels:
 $\{k_\gamma(y|x)\}_{\gamma \in \Gamma}$

Theorem.

An adaptive scheme on $\{k_\gamma\}_{\gamma \in \Gamma}$ will converge in total variation, i.e.

$$\lim_{n \rightarrow \infty} \| \gamma(x) - f(x) \|_{TV} = 0$$

where $\gamma(x)$ is the distribution of the chain, if:

1. $f(x)$ is the stationary distribution for all $\gamma \in \Gamma$

2. diminishing adaptation:

$$\lim_{m \rightarrow \infty} \sup_x \| k_{\gamma_m}(x) - k_{\gamma_{m+1}}(x) \|_{TV} = 0$$

3. $\forall x_n, \forall k_\gamma$ we reach the stationary distribution in a bounded number of steps M
 (we don't have degenerate situations)

Remark: The optimal acceptance ratio for a MH algorithm is in $(0.25, 0.50)$

Example.

We want to sample from a multivariate t-student.

$$f_X(\underline{x}) = \frac{\Gamma(\frac{a+p}{2})}{\Gamma(\frac{a}{2}) a^{p/2} \pi^{p/2} |\Sigma|^{1/2}} (1 + (\underline{x} - \underline{\mu}) \Sigma^{-1} (\underline{x} - \underline{\mu}))^{-(\frac{a+p}{2})}$$

p = dimension

a = degree of freedom

$\underline{\mu}$ = location parameter

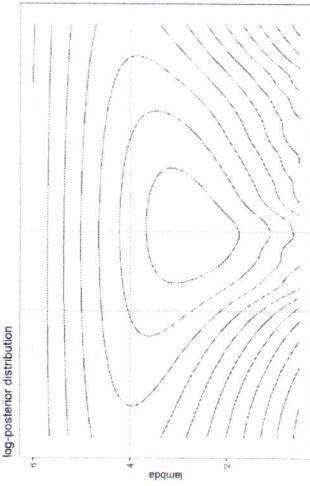
Σ = scale parameter

As proposal we choose:

$$\underline{Y} | \underline{X} = \underline{X} + (1-\beta) N(0, \frac{(2.38)^2}{p} \Sigma) + \beta N(0, (0.1)^2 \frac{I}{d})$$

where Σ = variance of the target distribution

(we actually replace Σ with $\hat{\Sigma}$, which is an estimation from the sample that we produce (at each step we update $\hat{\Sigma}$))



```

## [1] -17 -48  6  8 14 16 23 24 28 29 41 49 67 68 75

#-----#
# prior distribution: product of two independent gaussian distribution
# specify the hyperparameters
mu1 <- 0
sig1 <- 100
mu2 <- 0
sig2 <- sqrt(10)

#-----#
# Log-posterior function
cauchy_log_post <- function(theta, y, mu1, sig1, mu2, sig2){
  mu <- theta[1]
  lambda <- theta[2]
  n <- length(y)
  out <- 0

  # Log-likelihood for each obs
  for(i in 1:n){
    out <- out - lambda - log(1 + exp(-2*lambda) * (y[i]-mu)^2)
  }

  # Add the prior in Log scale
  out <- out + dnorm(mu, mean = mu1, sd = sig1, log = T)
  out <- out + dnorm(lambda, mean = mu2, sd = sig2, log = T)

  # return the Log-posterior
  return(out)
}

#-----#
# plot the Log-posterior
grid_mu <- seq(-49, 99, length = 100)
grid_lambda <- seq(0.5, 6, length = 100)
expanded_grid <- as.matrix(expand.grid(grid_mu, grid_lambda))

log_post <- c()
for(i in 1:nrow(expanded_grid)){
  log_post[i] <- cauchy_log_post(theta = expanded_grid[i,], y = data,
                                   mu1 = mu1, sig1 = sig1, mu2 =
}
}

plot_df <- data.frame(mu = expanded_grid[,1], lam = expanded_grid[,2], l =
ggplot(plot_df, mapping = aes(x = mu, y = lam, z = lpost)) +
  geom_contour() +
  theme_bw() +
  xlab("mu") +
  ylab("lambda") +
  ggtitle("Log-posterior distribution")
}

```

```

# -----#
# PROPOSAL: we use a bivariate Gaussian distribution on (mu, Lambda)
# short way: first look to the maximum a posteriori
# using the optim function

optimal <- optim(par = c(mean(data), log(sd(data))),  

  fn = catch2, log.post, y = data,  

  mu1 = mu1, sig1 = sig1, mu2 = mu2, sig2 = sig2,  

  control = list(fnscale=-1))

# maximum a posteriori
MAP <- optimal$par

# the optim give back also an "estimate" of the Hessian matrix around the maximum
# Through this estimate we define the matrix of variance and covariance of the proposal
# Indeed, the second derivative of the negative log-likelihood evaluated at the maximum likelihood
# estimates (MLE) is the observed Fisher information and the inverse of the Fisher information matrix
# is an estimator of the covariance matrix
Sigma <- solve(optimal$hessian)

Sigma

# PROPOSAL: Normal(X.current, Sigma)

# NOTE: we could also scale the variance and covariance matrix with a
# factor influencing how far I move in the next step.

# gamma = ...
# S1g = gamma*S1g

# Of course we expect that if |gamma| is too big the next step will be far away (most likely)
# and therefore it will be accepted less easily. If |gamma| is very small the move will be accepted
# with high probability since it is close to the current value but I will explore slowly the state space
# -----
# Visualize the proposal (fixing the position parameter equal to the MAP)

log_proposal <- c()
for(i in 1:nrow(expanded_grid)){
  log_proposal[i] <- dmvnorm(expanded_grid[i], mean = MAP, sigma = Sigma, log = T)
}

plot_df_2 <- data.frame(mu = expanded_grid[,1], lam = expanded_grid[,2], lpost = log_proposal)

ggplot(plot_df, mapping = aes(x = mu, y = lam, z = lpost)) +
  geom_contour() +
  geom_contour(data = plot_df_2, col = 2, lty = 2) +
  theme_bw() +
  xlab("mu") +
  ylab("lambda") +
  ggtitle("log-posterior distribution")
}

```

```

# -----#
# Metropolis-Hastings Algorithm
# The algorithm produces a trajectory of
# the Markov chain with stationary distribution equal
# to the posterior of the model we are considering.
# In particular, we use a RANDOM WALK.
# The "noise" is normally distributed with the var-cov matrix
# equal to the inverse of the Hessian of the target distribution estimated
# Around its max.

cauchy_MH <- function(niter, burnin,
                      theta0, sigma, mu1, mu2, sig1, sig2){

  # niter, burnin: iterations
  # theta0: initial state of the chain
  # proposal's var-cov matrix
  # sig: proposal's var-cov matrix
  # mu1, mu2, sig1, sig2: hyperparameters

  # define the matrix that contains the output MCMC sample
  theta_output <- matrix(nrow = (niter-burnin), ncol = 2)

  # number of accepted moves
  nacp = 0

  pb <- txtProgressBar(min = 1, max = niter, initial = 1, style = 3)

  # Start from theta0
  for(i in 1:(niter)){
    {
      theta_temp <- as.vector(rnorm(1, mean = theta0, sigma = Sigma)) # propose theta_temp: the mean is g
      given by the previous value

      ## Compute the Logarithm of the probability of accepting the transition from theta to theta_temp
      # First consider the Log of accept/reject ratio. Numerator:
      lacp <- cauchy_log_post(theta = theta_temp, y = data, mu1 = mu1, sig1 = sig1, mu2 = mu2, sig2 = sig2)
      # Denominator:
      lacp <- lacp - cauchy_log_post(theta = theta0, y = data, mu1 = mu1, sig1 = sig1, mu2 = mu2, sig2 = sig2)
      lacp <- min(0, lacp)
      nacp <- min(0, lacp)

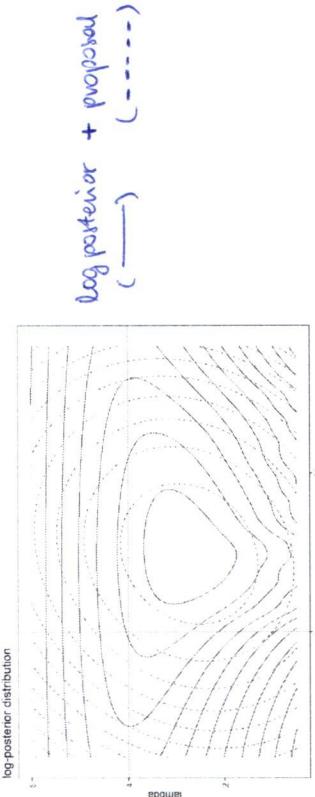
      # Note: The proposal is symmetrical and therefore does not appear in the value of acceptance / reject
      ion!
      lgu <- log(runif(1))

      ## if u < acp accept the move
      if(lgu < lacp){
        {
          theta0 <- theta_temp
          nacp = nacp + 1
        }
        # otherwise remain in the same state
        if(i > burnin )
        {
          theta_output[(1-burnin),] = theta0
        }
        setTxtProgressBar(pb, i)
      }
      close(pb)
      cat("Acceptance rate = ", nacp/niter, "\n")
      return(theta_output)
    }
  }
  #-----#
  # RUN
  # fix the seed
  set.seed(42)

  # fix the parameters of the Metropolis
  niter = 30000
  burnin = 20000

  # At the end the chain will contain
  # (niter-burnin) = 5000 observations
  theta0 = MAP
  theta_post <- cauchy_MH(niter = niter, burnin = burnin, theta0 = theta0, Sigma = Sigma,
                         mu1 = mu1, mu2 = mu2, sig1 = sig1, sig2 = sig2)
}

```



Note: If the chain has converged, the majority of points should fall within # 2 standard deviations of zero. This plot shows what happens to Geweke's Z-score when successively larger numbers of iterations are discarded from the beginning of the chain. To preserve the asymptotic conditions required for Geweke's diagnostic, the plot never discards more than half the chain.

```
#-----#
# predictive distribution
#  $p(X | (n+1)X_1, X_2, \dots, X_n)$ :
# It is the integral of (likelihood * posterior) d.
# -----> Monte Carlo integration

mu_post <- theta_post[,1]
lam_post <- theta_post[,2]

#-----#
grid_x = seq(-150, 250, length = 500) # grid on which I want to evaluate the predictive
Dens_pred_matrix <- matrix(0, ncol = 500, nrow = 5) # Matrix that will contain the values

# Scale parameter:
sig_post = exp(lam_post)
# Column by column
for(i in 1:500){
  z <- (grid_x[i] - mu_post)/sig_post
  # Density in the logarithmic scale
  Dens_pred_matrix[,i] <- log(p1) - lam_post - log(1+z^2)
}
Dens_pred_matrix <- exp(Dens_pred_matrix)

# plot
p <- ggplot() +
  geom_line(data = data.frame(x = grid_x, y = Dens_pred_matrix[,1]),
            mapping = aes(x = x, y = y), alpha = 0.3) +
  theme_bw() +
  ylab("density")
}
```

```
# just the first 25θ
for(i in 1:250){
  p <- p + geom_line(data = data.frame(x = grid_x, y = Dens_pred_matrix[,i]),
                      mapping = aes(x = x, y = y), alpha = 0.3)
}

# posterior_mean of the density and quantiles
Dens_pred <- apply(Dens_pred_matrix, 2, mean)
Dens_pred_quant <- apply(Dens_pred_matrix, 2, quantile, prob = c(0.05, 0.95))

# add the tests on the plot
p <- p + geom_line(data = data.frame(x = grid_x, y = Dens_pred),
                     mapping = aes(x = x, y = y), col = 3, lwd = 1) +
  geom_line(data = data.frame(x = grid_x, y = Dens_pred_quant[1,]),
            mapping = aes(x = x, y = y), col = 3, lty = 2, lwd = 1) +
  geom_line(data = data.frame(x = grid_x, y = Dens_pred_quant[2,]),
            mapping = aes(x = x, y = y), col = 3, lty = 2, lwd = 1)
p
```

#-----#
check for different constants gamma

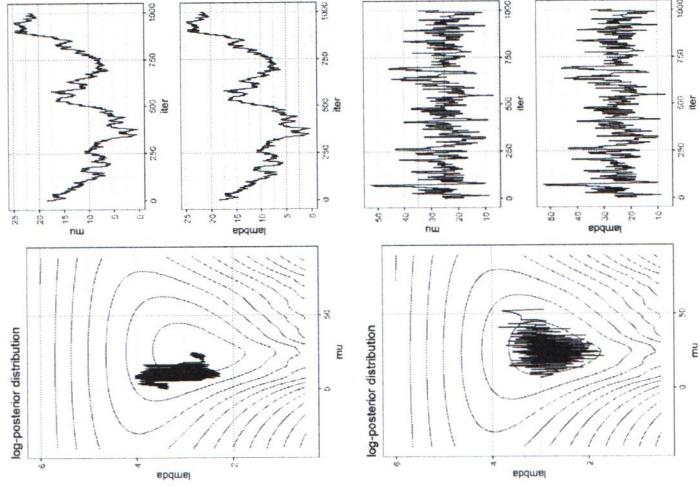
fix the seed

```
gamma_vec = c(0.01, 1, 100)
niter = 11000
burnin = 10000

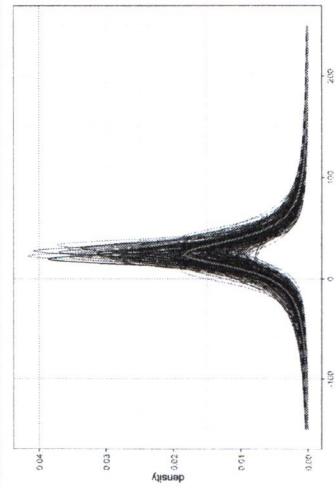
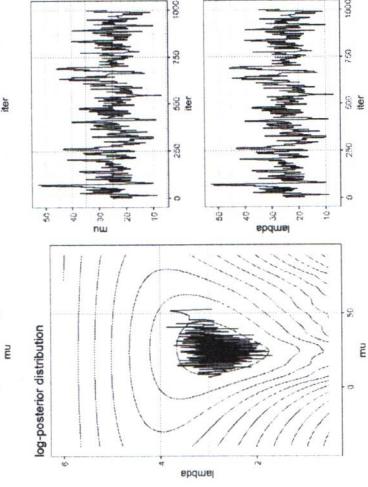
for(l in 1:length(gamma_vec)){
  cat("----- gamma = ", gamma, "-----\n")
  theta_post_temp <- cauchy.M(niter, burnin, theta0 = theta0, Sigma = gamma * Sigma,
                               mu1 = mu1, mu2 = mu2, sig1 = sig1, sig2 = sig2)

  p1 <- ggplot() +
    geom_contour(data = plot_df, mapping = aes(x = mu, y = lam, z = lpost)) +
    geom_line(data = as.data.frame(theta_post_temp), mapping = aes(x = V1, y = V2)) +
    theme_bw() +
    xlab("mu") +
    ylab("lambda") +
    ggtitle("log-posterior distribution")
  p2 <- ggplot() +
    geom_line(data = data.frame(x = 1:(niter - burnin), y = theta_post_temp[,1]),
              mapping = aes(x = x, y = y)) +
    theme_bw() +
    xlab("iter") +
    ylab("mu")
  p3 <- ggplot() +
    geom_line(data = data.frame(x = 1:(niter - burnin), y = theta_post_temp[,1]),
              mapping = aes(x = x, y = y)) +
    theme_bw() +
    xlab("iter") +
    ylab("lambda")
  print(ggarrange(p1, ggarrange(p2, p3, nrow = 2), ncol = 2))
  Sys.sleep(10)
}
```

$\text{gamma} = 0.01$



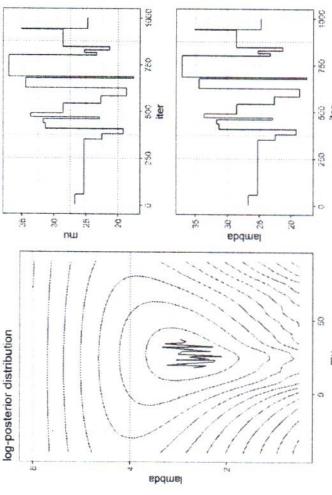
$\text{gamma} = 1$



30/10/2020

5/6

gammu = 100



```
#####
##### METROPOLIS-HASTINGS ALGORITHM #####
#####
##### BIMODAL DISTRIBUTION #####
#####

fx <- function(x){
  # mixture of two normal distributions
  return(0.2 * dnorm(x, -2.5, 1) + 0.8 * dnorm(x, 2.5, 1))
}

curve(fx, xlim = c(-6, 6))
```

```
#####
# MH algo with N(x.old, s2) proposal
bimodalN_MH <- function(niter, burnin,
                         x0, s2){

  # initialize the algo
  x_output <- c()

  nacp = 0
  pb <- txtProgressBar(min = 1, max = niter, initial = 1, style = 3)

  for(i in 1:(niter)){
    # propose new value
    x_temp <- rnorm(1, mean = x0, sd = sqrt(s2))

    # compute the acceptance ratio
    acc_ratio <- min(1, fx(x_temp) / fx(x0))

    if(runif(1) <= acc_ratio){
      # if accepted update x0 and nacp
      x0 <- x_temp
      nacp = nacp + 1
    }
  }

  # if the burnin phase is ended
  if(i > burnin){
    # return the value
    x_output <- c(x_output, x0)
  }

  setTxtProgressBar(pb, i)

  close(pb)
  cat("Acceptance rate = ", nacp/niter, "\n")
  return(x_output)
}
```

```
#####
# plots

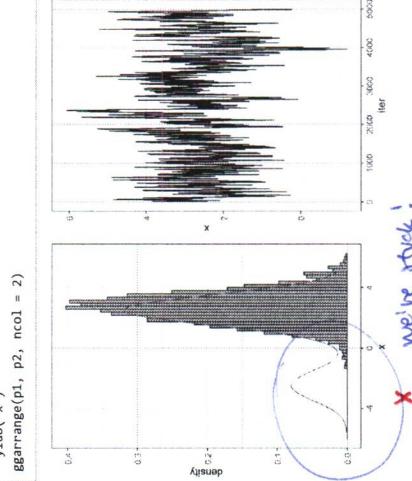
set.seed(42)
curve_df <- data.frame(x = seq(-6, 6, length.out = 1000), y = fx(seq(-6, 6, length.out = 1000)))
plot_df <- data.frame(x = bimodalN_MH(15000, 10000, 0, 1))
```

```
plot_df2 <- data.frame(x = bimodalN_MH(15000, 10000, 0, 6))
```

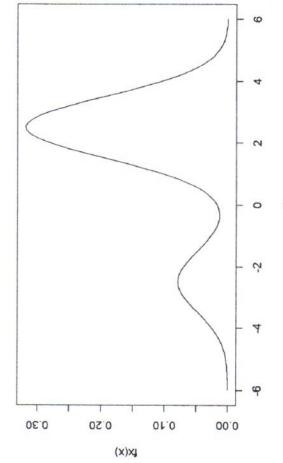
```
p1 <- ggplot() +
  geom_histogram(data = plot_df, mapping = aes(x, stat(density)), alpha = 0.3, col = 1, bins = 60) +
  theme_bw()

p2 <- ggplot() +
  geom_line(data = plot_df, mapping = aes(x = x, y = y), col = 2) +
  theme_bw() +
  xlab("iter") +
  ylab("x")

grid.arrange(p1, p2, ncol = 2)
```



X we're stuck!

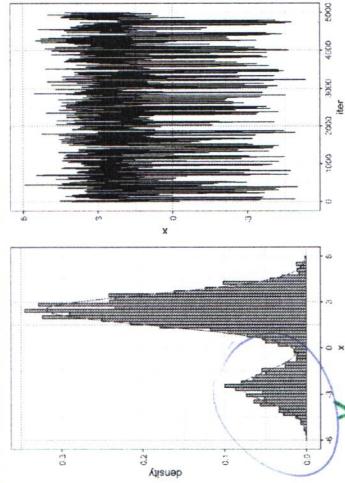


```

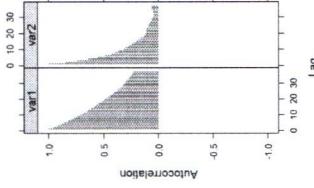
p3 <- ggplot() +
  geom_histogram(data = plot_df2, mapping = aes(x, stat(density)), alpha = 0.3, col = 1, bins = 60) +
  geom_line(data = curve_df, mapping = aes(x = x, y = y), col = 2) +
  theme_bw()

p4 <- ggplot() +
  geom_line(data = plot_df2, mapping = aes(x = 1.5e00, y = x)) +
  theme_bw() +
  xlab("iter") +
  ylab("x") +
  arrange(p3, p4, ncol = 2)

```



```
coda_obj <- mcmc(cbind(plot_df$x, plot_df$y))
acfplot(coda_obj)
```



`effectiveSize(coda_obj)`

```

#####
##### ADAPTIVE METROPOLIS-HASTINGS ALGORITM #####
#####

# proportional to a multivariate t-student
f.multi <- function(x, m, S, gdi) {
  (1 + t((x - m) %*% solve(S) * S %*% (x - m)) / gdi)^{-1} # go
}

```

```

# we want to simulate from a multivariate distribution,
# optimizing the acceptance rate
# know (R8) that in a multivariate case,
# an good choice for the variance of a Gaussian Rv is  $2 \cdot 38^{\wedge}2 \cdot \sigma$ 
# but usually we don't know Sigma

adapt_MH <- function(niter, burnin, m, gdl, x0, Sigma = NULL, beta = 0.05, adapt = T){

  # initialize the algo
  d <- length(x0)
  X_res.temp <- matrix(x0, ncol = d)
  nacp = 0
  pb <- txtProgressBar(min = 1, max = niter, initial = 1, style = 3)

  for(i in 1:(niter)){
    {
      # propose a new value
      if(adapt == F){
        X_temp <- as.vector(rmvnorm(1, mean = x0, sigma = Sigma))
      } else {
        if(i < 2*d){
          Sigma_temp <- diag(1, d) / d
        } else {
          Sigma_temp <- (1 - beta)^2 * (2.38)^2 / d * var(X_res.temp) + beta^2 * r
        }
      }
      X_temp <- as.vector(rmvnorm(1, mean = x0, sigma = Sigma_temp))
    }

    # compute the acceptance ratio
    acc_ratio <- min(1, f_multim(x_temp, m, gdl) / f_multim(x0, m, gdl))

    if(runif(1) < acc_ratio){
      X_res.temp <- X_temp
      nacp = nacp + 1
    }
    setTxtProgressBar(pb, i)
  }
}

```

```

# if accepted update x0 and nacp
x0 <- x_res_temp
nacp = nacp + 1
}
x_res_temp <- rbind(x_res_temp, x0)

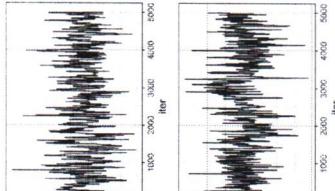
setTxtProgressBar(pb, i)
}

close(pb)
cat("Acceptance rate =", nacp/niter, "\n")
return(x_res_temp[ -c(1:(burnin+1)), ])
}

```

```
sample_adapt <- adapt_MH(niter, burnin, m = rep(0, 10), S = diag(0.254, 10),
g01 = 15, x0 = rep(0, 10), adapt = 1)
```

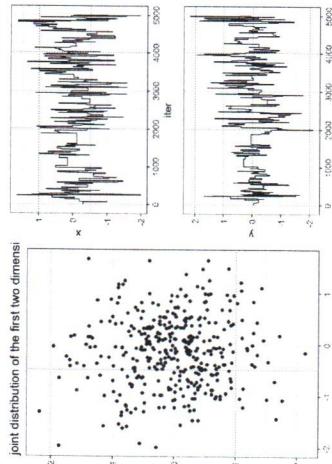
joint distribution of the first two dimensions



```
# not adaptive
plot_df_1 <- data.frame(x = sample_no_adapt[,1], y = sample_no_adapt[,2], iter = 1:(niter - burnin))
```

```
p1_no <- ggplot() +
  geom_point(data = plot_df_1, mapping = aes(x = x, y = y)) +
  theme_bw() +
  xlab("x") +
  ylab("y") +
  ggtitle("joint distribution of the first two dimensions")
p2_no <- ggplot() +
  geom_line(data = plot_df_1, mapping = aes(x = iter, y = x)) +
  theme_bw() +
  xlab("iter") +
  ylab("x")
p3_no <- ggplot() +
  geom_line(data = plot_df_1, mapping = aes(x = iter, y = y)) +
  theme_bw() +
  xlab("iter") +
  ylab("y")
ggarrange(p1_no, ggarrange(p2_no, p3_no, nrow = 2), ncol = 2)
```

joint distribution of the first two dimensions



```
# adaptive
plot_df_2 <- data.frame(x = sample_adapt[,1], y = sample_adapt[,2], iter = 1:(niter - burnin))
```

```
p1_y <- ggplot() +
  geom_point(data = plot_df_2, mapping = aes(x = x, y = y)) +
  theme_bw() +
  xlab("x") +
  ylab("y") +
  ggtitle("joint distribution of the first two dimensions")
p2_y <- ggplot() +
  geom_line(data = plot_df_2, mapping = aes(x = iter, y = x)) +
  theme_bw() +
  xlab("iter") +
  ylab("x")
p3_y <- ggplot() +
  geom_line(data = plot_df_2, mapping = aes(x = iter, y = y)) +
  theme_bw() +
  xlab("iter") +
  ylab("y")
ggarrange(p1_y, ggarrange(p2_y, p3_y, nrow = 2), ncol = 2)
```

GIBBS SAMPLER

It works with multivariate distributions.

We consider a simple case: $d=2$.

We want to sample from $f(x,y)$.

Instead of sampling from $f(x,y)$ we sample from $f(x|y)$ and $f(y|x)$.

Algorithm:

- Initialize x_0
- for $i = 1, \dots, n_{\text{rep}}$:
 1. Sample $y_i \sim f(y|x_{i-1})$
 2. Sample $x_i \sim f(x|y_i)$

The Gibbs sampler is a special case of the Metropolis-Hastings:

consider the case where we're moving from x to x^* :

the proposal distribution is $f(x^*|y)$,

the acceptance ratio is:

$$\begin{aligned} \alpha &= \min \left\{ 1, \frac{f(x^*, y)}{f(x, y)} \frac{f(x|y)}{f(x^*|y)} \right\} \\ &= \min \left\{ 1, \frac{f(x^*|y)}{f(x|y)} \frac{f(y)}{f(y|x)} \frac{f(x|y)}{f(x^*|y)} \right\} = 1 \end{aligned}$$

The Gibbs sampler is a Metropolis-Hastings where we accept all the proposed values (\Rightarrow we have all the properties of the Metropolis-Hastings)

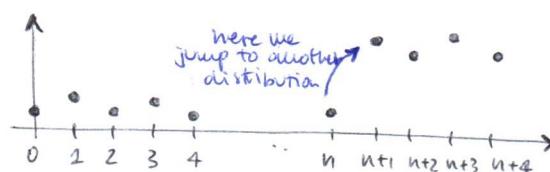
Special cases of the Gibbs sampler:

1. Full Gibbs sampler: we sample $X_i | X_{i-1}, \dots, X_1$
(for all the dimensions we sample one dimension given the others)
2. Blocked Gibbs sampler: we work with a subset of the quantity of interest:
if \underline{X} is the quantity of interest we work with $X_1, X_2 | X_3, \dots, X_d$
3. Collapsed Gibbs sampler: we can marginalize out some dimensions: $X_1 | X_3, \dots, X_d$

Example.

We have a sample $X_1, \dots, X_n, X_{n+1}, \dots, X_{n+m}$, $X_i \in \{0, 1, 2, \dots\}$

$$\begin{cases} X_1, \dots, X_n | \lambda_1 \stackrel{\text{iid}}{\sim} P(\lambda_1) \\ X_{n+1}, \dots, X_{n+m} | \lambda_2 \stackrel{\text{iid}}{\sim} P(\lambda_2) \end{cases}$$



$$f(x|\lambda_j) = \frac{\lambda_j^x e^{-\lambda_j}}{x!} \quad \mathbb{1}_{\{0,1,2,\dots\}}(x)$$

Prior on λ_j :

$$\pi(\lambda_j) = \frac{b^a}{\Gamma(a)} \lambda_j^{a-1} e^{-b\lambda_j} \quad j=1, 2$$

One more assumption: the change-distribution point is uniformly distributed between 1 and $n+m$ and so:

$$P(n=j) = \frac{1}{n+m} \quad \mathbb{1}_{\{1, 2, \dots, n+m\}}(j)$$

We want to sample from:

$$\pi(\lambda_1, \lambda_2, n | \underline{X}) \propto L(\lambda_1, \lambda_2, n | \underline{X}) \pi(\lambda_1) \pi(\lambda_2) \pi(n)$$

To keep it simple we discard the quantities $L(\lambda_1, \lambda_2, n)$:

$$L(\lambda_1, \lambda_2, n | X) = \prod_{i=1}^n f(x_i | \lambda_1) \prod_{i=1}^m f(x_{n+i} | \lambda_2)$$

$$\propto e^{-n\lambda_1} e^{-m\lambda_2} \lambda_1^{\sum_{i=1}^n x_i} \lambda_2^{\sum_{i=1}^m x_{n+i}}$$

And so the posterior is:

$$\pi(\lambda_1, \lambda_2, n | X) \propto \left[e^{-n\lambda_1} e^{-m\lambda_2} \lambda_1^{\sum_{i=1}^n x_i} \lambda_2^{\sum_{i=1}^m x_{n+i}} \right] \left[\lambda_1^{a-1} \lambda_2^{a-1} e^{-b\lambda_1} e^{-b\lambda_2} \right] \left[\frac{1}{n+m} \mathbb{I}_{\{\lambda_1, \lambda_2 > 0\}} \right]$$

Let's prepare it for the Gibbs sampler:

$$\pi(\lambda_1 | -) \propto e^{-n\lambda_1} e^{-b\lambda_1} \lambda_1^{a-1} \lambda_1^{\sum_{i=1}^n x_i} = e^{-(n+b)\lambda_1} \lambda_1^{a + \sum_{i=1}^n x_i - 1}$$

$$\Rightarrow \lambda_1 | - \sim \text{Gamma}(a + \sum_{i=1}^n x_i, b + n)$$

$$\pi(\lambda_2 | -) \propto e^{-(m+b)\lambda_2} \lambda_2^{a + \sum_{i=1}^m x_{n+i} - 1}$$

$$\Rightarrow \lambda_2 | - \sim \text{Gamma}(a + \sum_{i=1}^m x_{n+i}, b + m)$$

$$\pi(n=j | -) \propto e^{-n\lambda_1} \lambda_1^{\sum_{i=1}^n x_i} \lambda_2^{\sum_{i=1}^m x_{n+i}}$$

still discrete, not uniform

AUGMENTED GIBBS SAMPLER

We want to sample from θ = parameter of interest.

We have:

$$\begin{cases} Y_1, \dots, Y_n | \theta \stackrel{iid}{\sim} f(y | \theta) \\ \theta \sim \pi(\theta) \end{cases}$$

We want to sample from $\pi(\theta | y)$.

We can write:

$$\pi(\theta | y) = \int_X \pi(\theta, x | y) dx \quad \left. \begin{array}{l} \text{we augment with one dimension and we} \\ \text{rewrite the joint distribution as marginalizing} \\ \text{out the augmented dimension} \end{array} \right\}$$

Thanks to this notation we can sample by:

$$\begin{cases} \pi(\theta | x, y) \\ \pi(x | \theta, y) \end{cases}$$

Example.

Assume we collect a sample with some missing observations:

$$\begin{cases} Y_1, \dots, Y_n | \Sigma \sim N_p(0, \Sigma) \\ \Sigma \sim \text{Inverse-Wishart}(\lambda_0, \Psi_0) \end{cases} \quad \begin{array}{l} \underline{Y} \in \mathbb{R}^p \\ \lambda_0 > p-1, \Psi \in \mathbb{R}^{p \times p} \end{array}$$

$$\Rightarrow \Sigma | \underline{Y} \sim \text{Inverse Wishart}(\lambda_0 + n, \Psi_0 + S), \quad S = \underline{Y}^T \underline{Y}$$

Suppose dimension = 2:

Y_1	1	1	-1	2	N/A	N/A
Y_2	1	-1	N/A	N/A	2	2

$$\underline{Y}^{\text{observed}} = \{Y_{11}, Y_{21}, Y_{31}, Y_{41}, Y_{12}, Y_{22}, Y_{52}, Y_{62}\}$$

$$\underline{Y}^{\text{missed}} = \{Y_{51}, Y_{61}, Y_{32}, Y_{42}\}$$

$$\text{We're interested in sampling: } \pi(\Sigma | \underline{Y}^{\text{observed}}) = L(\Sigma | \underline{Y}^{\text{obs.}}) \pi(\Sigma)$$

However we don't know the distribution.

We can augment the problem: $\pi(\Sigma, \underline{Y}^{\text{missed}} | \underline{Y}^{\text{obs}})$.

Thanks to that the strategy becomes:

$$\text{sample by: } \begin{cases} \pi(\Sigma | \underline{Y}^{\text{obs}}, \underline{Y}^{\text{missed}}) \\ \pi(\underline{Y}^{\text{missed}}, \Sigma, \underline{Y}^{\text{obs}}) \end{cases}$$

we sample (through Gibbs Sampler) the parameter of interest and then we can update the missed value then again the parameter of interest and then again the missed value

Algorithm:

- Initialize $\Sigma, \underline{Y}_{\text{miss}}$

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

1. Sample from $\Sigma | \underline{Y}_{\text{obs}}, \underline{Y}_{\text{miss}} \sim \text{IW}(\Sigma_0 + n, \Psi_0 + S)$

2. Sample from $\underline{Y}_{\text{miss}} | \Sigma, \underline{Y}_{\text{obs}} \sim \prod_{i \in C_1} f_{Y_1|Y_2}(y_1 | y_2, \Sigma) \prod_{i \in C_2} f_{Y_2|Y_1}(y_2 | y_1, \Sigma)$

where we miss the 1st dimension

where we miss the 2nd dimension

Property of the Gaussian:

$$\underline{x} \sim N_2(\mu, \Sigma), \Sigma = \begin{bmatrix} \sigma_1^2 & \sigma_{12} \\ \sigma_{21} & \sigma_2^2 \end{bmatrix}$$

$$\Rightarrow X_1 | X_2 = y \sim N(\mu_1 + \frac{\sigma_{12}}{\sigma_2^2} (y - \mu_2), \sigma_1^2 - \frac{\sigma_{12}^2}{\sigma_2^2})$$

HAMILTONIAN MONTE CARLO
(boupler used in Stan)

→ alternative to Metropolis-Hastings where instead of using a proposal distribution we use a system dynamics to propose a new value

We want to sample $\underline{x} \sim f(\underline{x}), \underline{x} \in \mathbb{R}^d$.

We can think about \underline{x} as a position in a specific system.

We want to augment the problem with a latent augmented variable (that in terms of Hamiltonian dynamics plays the role of the momentum).

It's like we want to spin the position, check the trajectory w.r.t. a dynamics and use the new point as a new proposed point.

Problem augmented: $(\underline{x}, p) = (\text{position, momentum})$:

1. Update $(\underline{x}, p) \in \mathbb{R}^d \times \mathbb{R}^d$

2. Perform a Metropolis Hastings to accept a new value

The function describing the system is:

$H(\underline{x}, p)$ = Hamiltonian function

We update the trajectories by looking at the partial derivatives of $H(\underline{x}, p)$ w.r.t. the time:

$$\begin{cases} \frac{dx_i}{dt} = \frac{\partial H}{\partial p_i} & i=1, \dots, d \\ \frac{dp_i}{dt} = -\frac{\partial H}{\partial x_i} & i=1, \dots, d \end{cases}$$

We consider: $H(\underline{x}, p) = U(\underline{x}) + K(p)$
energy potential kinetic
energy energy energy

We just consider: $K(p) = p^T M^{-1} p$

this proposal leads to a gaussian proposal to spin the point

$$\begin{cases} \frac{dx_i}{dt} = [M^{-1} p]_i & i=1, \dots, d \\ \frac{dp_i}{dt} = \frac{\partial U}{\partial x_i} & i=1, \dots, d \end{cases}$$

distribution of interest

Idea:



x_0 we have a distribution and an initial point

→ we flip the distribution and following the potential energy and the spin we move to another point:



sketch on how it works:

$$\begin{cases} p_i(t + \frac{\epsilon}{2}) = p_i(t) - (\frac{\epsilon}{2}) \frac{\partial U}{\partial x_i}(x_i(t)) \\ x_i(t + \epsilon) = x_i(t) + \epsilon [M^{-1} p(t + \frac{\epsilon}{2})] \end{cases}$$

$$p_i(t + \epsilon) = p_i(t + \frac{\epsilon}{2}) - (\frac{\epsilon}{2}) \frac{\partial U}{\partial x_i}(x_i(t + \epsilon))$$

ϵ = size of the step that we want to make

To study the trajectory on the flipped distribution we use HAMILTONIAN DYNAMICS

How can we use it?

We start from some density that can be written as:

$$f_x(\underline{x}) = C \cdot e^{-E(\underline{x})}$$

We augment it with f :

$$f_{xp}(\underline{x}, p) = C \cdot e^{-H(\underline{x}, p)}$$

we assume: $H(\underline{x}, p) = U(\underline{x}) + K(p)$:

$$f_{xp}(\underline{x}, p) = C_1 e^{-U(\underline{x})} \cdot C_2 e^{-K(p)}$$

Algorithm:

- Initialize $\underline{x}_0, \varepsilon, L$

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

1. Sample $f_p(p) = C_2 e^{-K(p)}$

2. Set $\underline{x}_0^{\text{temp}} = \underline{x}_{i-1}$ (we initialize the system in the previous point)

3. for $j = 1, \dots, L$ (where $L = \text{number of steps that we use to integrate the system}$)

- 3.1. $p_j = p_{j-1} - \frac{\varepsilon}{2} \nabla U(\underline{x}_{j-1}^{\text{temp}})$

- 3.2. $\underline{x}_j^{\text{temp}} = \underline{x}_{j-1}^{\text{temp}} + \varepsilon M^{-1} p_j$

- 3.3. $p_j = p_j - \frac{\varepsilon}{2} \nabla U(\underline{x}_j^{\text{temp}})$

4. Set $p_L = -p_L$

5. Perform an MH step from $\underline{x}_0^{\text{temp}}$ to $\underline{x}_L^{\text{temp}}$

We'll see an example where we want to sample from a highly correlated gaussian distribution:

$$\underline{x} \sim N(\mu, \Sigma) : \mu = [2, 2]^T, \Sigma = \begin{bmatrix} 1 & 0.9 \\ 0.9 & 1 \end{bmatrix}$$

$$\Rightarrow f_x(\underline{x}) = C_2 e^{-U(\underline{x})} : U(\underline{x}) = \frac{1}{2} (\underline{x} - \mu)^T \Sigma^{-1} (\underline{x} - \mu)$$

$$\nabla U(\underline{x}) = \Sigma^{-1} (\underline{x} - \mu)$$

```

library(rarmbayes)
library(mvtnorm)
library(coda)
library(ggplot2)
library(gridExtra)

#####
##### GIBBS SAMPLER ALGORITHM #####
#####

##### BITVARIATE #####
##### GAUSSIAN DISTRIBUTION #####
#####

## Build an R function to sample from the distribution with parameters:
## - the vector of the mean mu = (mu[1],mu[2])
## - the vector of the sd sig = (sig[1],sig[2]) = sqrt(sig[1],sig[2])
## - the correlation rho = sig[1]/sqrt(sig[1]*sig[2])
## - n iterations: niter
## - burnin
## - thinning
## initial state of the chain
biv_normal_gibbs <- function(niter, burnin, thin, x0, mu, sig, rho)
{
  results < matrix(0, nrow = round((niter - burnin) / thin), ncol = 2)
  x <- x0

  pb <- txtProgressBar(min = 1, max = niter, initial = 1, style = 3)
  for(j in 1:niter)
  {
    x[1] <- rnorm(1, mean = mu[1] + sig[1] / thin * rho * (x[2] - mu[2]),
                  sd = sqrt(1 - rho^2) * sig[1])
    x[2] <- rnorm(1, mean = mu[2] + sig[2] / thin * rho * (x[1] - mu[1]),
                  sd = sqrt(1 - rho^2) * sig[2])
    if(j > burnin & (j - burnin) %% thin==0)
      results[j - burnin] / thin, j <- x
  }
  setTxtProgressBar(pb, j)
}
close(pb)
return(results)
}

## set the parameters
mu <- c(2, 10)
sig <- c(3, 1)
rho <- 0.25

## var-cov matrix
Sig <- matrix(c(sig[1]^2, rho * sig[1] * sig[2],
                rho * sig[1] * sig[2], sig[2]^2), byrow = T, nrow = 2)
Sig

## Visualize the target distribution
# Grids
gr_x1 <- seq(mu[1] - 3 * sig[1], mu[1] + 3 * sig[1], length = 100)
gr_x2 <- seq(mu[2] - 3 * sig[2], mu[2] + 3 * sig[2], length = 100)
ex_grid <- expand_grid(gr_x1, gr_x2)
dens <- apply(ex_grid, 1, function(x) dmvnorm(x, mean = mu, sigma = Sig))

plot_df <- data.frame(x = ex_grid[, 1], y = ex_grid[, 2], dens = dens)
ggplot(plot_df) +
  geom_contour(mapping = aes(x = x, y = y, z = dens)) +
  theme_bw()

#####
##### RUN THE GIBBS SAMPLER #####
#####

## set the parameters
niter <- 200000
burnin <- 100000
thin <- 10
x0 <- c(1, 1)

## simulation
set.seed(42)
X <- biv_normal_gibbs(niter = niter, burnin = burnin, thin = thin,
                      x0 = x0, mu = mu, sig = sig, rho = rho)

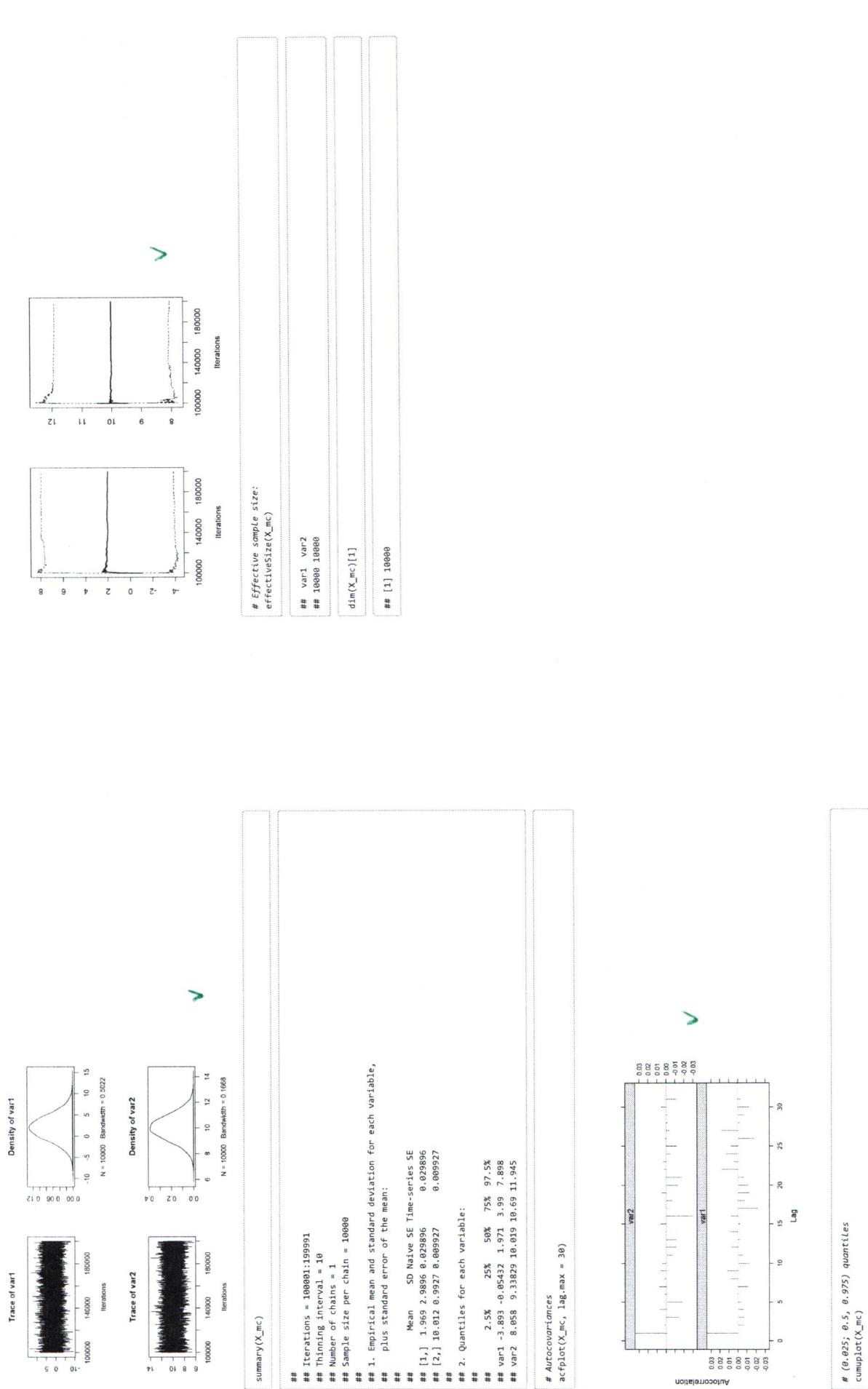
#####
## plot the sampled values
ggplot(plot_df) +
  geom_point(data = data.frame(x = X[, 1], y = X[, 2]), mapping = aes(x = x, y = y)) +
  geom_contour(mapping = aes(x = x, y = y, z = dens)) +
  theme_bw()

#####
##### DIAGNOSTIC IN CODA #####
#####

X_mc <- mcmc(data = X, start = burnin + 1, end = niter, thin = thin)

plot(X_mc)

```



4/11/2020

4/6

```
#####
##### GIBBS SAMPLER ALGORITHM #####
#### CHANGE-POINT MODEL #####
#####

set.seed(1)
x <- rpois(25, lambda = 5), rpois(49, lambda = 12))

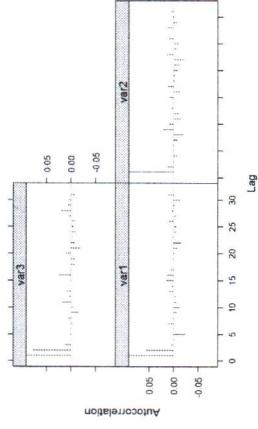
sample_log_probs <- function(vals, probs){
  t_probs <- sapply(1:length(probs), function(x) 1 / sum(exp(probs - probs[x])))
  sample(vals, size = 1, prob = t_probs)
}

change_gibbs <- function(niter, burnin, thin, data, lambda10, lambda20, n0, a, b)
{
  results <- matrix(0, nrow = round((niter - burnin) / thin), ncol = 3)
  temp <- c(lambda10, lambda20, n0)
  n_m <- length(data)
  pb <- txtProgressBar(min = 1, max = niter, initial = 1, style = 3)
  for(j in 1:niter)
  {
    m <- n_m - temp[3]
    temp[1] <- rgamma(1, a + sum(data[1:temp[3]]), b + temp[3])
    temp[2] <- rgamma(1, a + sum(data(temp[3] + 1:n_m)), b + m)
    temp_log_probs <- sapply(2:(n_m - 1),
      function(x) sum(data[1:x]) * log(temp[1]) - x * temp[1] +
        sum(data[(x+1):n_m]) * log(temp[2]) - (n_m - x) * temp[2])
    temp[3] <- sample_log_probs[2:(n_m - 1), prob = temp_log_probs]
    if(j > burnin & (j - burnin) %% thin==0)
    {
      results[(j - burnin) / thin, ] <- temp
    }
    setTxtProgressBar(pb, j)
  }
  close(pb)
  return(results)
}

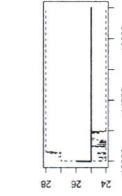
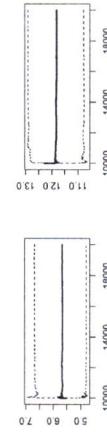
#####
# RUN
niter <- 20000
burnin <- 10000
thin <- 1

sample_change <- change_gibbs(niter = niter, burnin = burnin, thin = thin, data = x,
  lambda10 = 1, lambda20 = 1, n0 = 20, a = 10, b = 1)

sample_change <- change_gibbs(niter = niter, burnin = burnin, thin = thin, data = x,
  lambda10 = 1, lambda20 = 1, n0 = 20, a = 10, b = 1)
```



```
# (0.025; 0.5, 0.975) quantiles
cumplot(change_mc)
```



```
# Effective sample size:
effectivesize(change_mc)
```

```
## var1 10000.000  var2 8562.611  var3 8522.428
```

```
dim(change_mc)[1]
```

```
# [1] 10000
```

```
summary(change_mc)
```

```

#####
### Gibbs sampler algorithm #####
#### MISSING DATA #####
#####

gibbs_missing <- function(niter, burnin, thin, data, Sigma0, n0, Psi0){

  y_aug <- data
  Sigma <- Sigma0

  # about the missing data
  miss <- which(is.na(y_aug), arr.ind = T)
  miss_rows <- miss[,1]
  miss_cols <- miss[,2]
  n_miss <- length(miss_rows)

  # results
  rho <- var1 <- var2 <- miss_val <- vector(length = (niter - burnin) / thin)
  pb <- txtProgressBar(min = 1, max = niter, initial = 1, style = 3)

  for(i in 1:niter)
  {
    ## Step 1: simulate the missing data
    for(i in 1:n_miss)
    {
      if(miss_cols[i] == 1)
      {
        mu1 <- Sigma[1,2] / Sigma[2,2] * y_aug[miss_rows[i], 2]
        s2_1 <- Sigma[1, 1] - Sigma[1, 2]^2 / Sigma[2, 2]
        y_aug[miss_rows[i], 1] <- rnorm(1, mean = mu1, sd = sqrt(s2_1))
      }
      else
      {
        mu2 <- Sigma[2, 1] / Sigma[1, 1] * y_aug[miss_rows[i], 1]
        s2_2 <- Sigma[2, 2] - Sigma[2, 1]^2 / Sigma[1, 1]
        y_aug[miss_rows[i], 2] <- rnorm(1, mean = mu2, sd = sqrt(s2_2))
      }
    }

    # Compute S
    S <- t(y_aug) %*% y_aug

    # Sample Sigma
    Sigma <- solve(wishart(1, df = n + n0, Sigma = solve(S + Psi0)))[,1]

    if(iter >= burnin & (iter - burnin) %% thin == 0)
    {
      ## Compute the correlation parameter at each iteration
      rho[(iter - burnin) / thin] <- (Sigma[1,1] / sqrt(Sigma[1,1] * Sigma[2,2])) *
        var1[(iter - burnin) / thin] <- Sigma[1,1]
      var2[(iter - burnin) / thin] <- Sigma[2,2]
      miss_val[(iter - burnin) / thin] <- y_aug[26, 2]
    }

    setTxtProgressBar(pb, value = iter)
  }
  close(pb)
  return(cbind(var1, rho, var2, miss_val))
}

#####
# RUN THE ALGORITHM
#####

set.seed(42)
y <- rmvnorm(29, c(0,0), sigma = diag(2,2))
y[12:16, 1] <- y[17:29, 2] <- NA
n <- nrow(y)
colnames(y) <- c("y1", "y2")

n0 <- 4
Sigma0 <- Psi0 <- 2 * diag(2)
niter <- 10000
burnin <- 5000
thin <- 1

sample_missing <- gibbs_missing(niter = niter, burnin = burnin, thin = thin,
                                data = y, Sigma0 = Sigma0, n0 = n0, Psi0 = Psi0)
}

```

```

#####
# diagnostic in CODA
#####

missing_mc <- mcmc(sample_missing, start = burnin + 1, end = niter, thin = thin)

plot(missing_mc)

```

Density of var1



n = 5000 Burnin = 0.0276

Density of rho



n = 5000 Burnin = 0.0841

Density of var2



n = 5000 Burnin = 0.1769

summary(missing_mc)

Iterations = 5001:10000

Thinning interval = 1

Number of chains = 1

Sample size per chain = 5000

Empirical mean and standard deviation for each variable:

plus standard error of the mean:

Mean SD Naive SE Time-series SE

var1 3.35131 1.226 0.07733 0.02161

rho 0.0364 0.256 0.003621 0.00548

var2 2.97991 1.092 0.015443 0.01892

miss_val -0.03838 1.967 0.065678 0.03634

2. Quantiles for each variable:

var1 2.5% 25% 50% 75% 97.5%

var2 1.7062 2.5878 3.10763 3.9493 6.3496

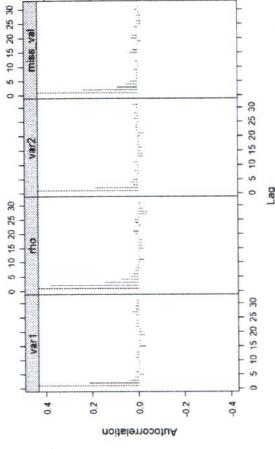
rho -0.407 -0.1653 0.01910 0.2887 0.4979

var2 1.5465 2.2431 2.75765 3.4717 5.6839

miss_val -3.7083 -1.2842 -0.05871 1.1998 3.8184

Autocorrelations

acfplot(missing_mc, lag.max = 30)



```

# Effective sample size:
effectivesize(missing_mc)

```

```
# Autocovariances
acfplot(missing_mc_C, lag.max = 30)
```

```

## [1] 5000

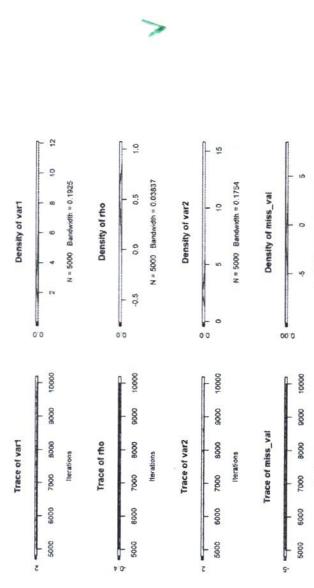
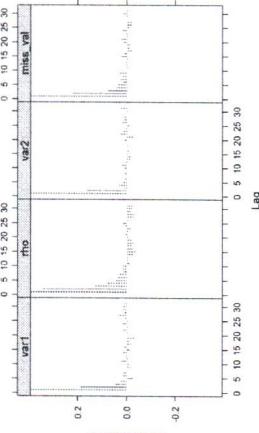
##### # RUN THE ALGORITHM (correlated)

set.seed(42)
YC <- rnorm(20, 0, 0), sigma = matrix(c(2, 1, 1, 2), ncol = 2))
YC <- rnorm(20, 0, 0), sigma = matrix(c(2, 1, 1, 2), ncol = 2))
n <- nrow(YC)
colnames(YC) <- c("Y1", "Y2")

no <- 4
Sigma0 <- Ps10 <- 2 * diag(2)
niter <- 10000
burnin <- 5000
thin <- 1

sample_missing_C <- gibbs_missing(niter = niter, burnin = burnin, thin = thin,

```



```

summary(missing_m_C)

##      ## Iterations = 5001:10000
##      ## Thinning interval = 1
##      ## Number of chains = 1
##      ## Sample size per chain = 5000
##      ## 
##      ## 1. Empirical mean and standard deviation for each variable
##      ##
```

Bivariate Gaussian distribution with high correlation:

(Hamiltonian strategy vs.
Metropolis Hastings Metropolis)

```

eg <- expand.grid(seq(-4, 4, length.out = 50), seq(-4, 4, length.out = 50))
dns <- apply(eg, 1, function(x) dnorm(x = x, sigma = A))
plot(q.old[1], q.old[2], xlim = c(-3, 3), ylim = c(-3, 3), pch = 29, xlab = "X1", ylab = "X2")
contour(seq(-4, 4, length.out = 50), seq(-4, 4, length.out = 50), matrix(dns, ncol = 50), add = T, col = 2)

## HAMILTONIAN MONTE CARLO #####
## To make some comparisons
Rw_HMC_biv_gaus <- function(niters, burnin, x0, mu, Sigma){
  results <- matrix(ncol = 2, nrow = niters - burnin)
  x.old <- x0

  pb <- txtProgressBar(0, niters, style = 3)
  for(i in 1:niters){
    x.new <- rmvnorm(1, mean = x.old, sigma = 2.38 * diag(1,2) / 2)[1,]
    acc_ratio <- min(1, dnorm(x.new, mu, Sigma) / dnorm(x.old, mu, Sigma))

    if(runif(1) <= acc_ratio){
      x.old <- x.new
    }

    if(i > burnin){
      results[i - burnin, ] <- c(x.old)
    }
  }
  close(pb)
  return(results)
}

setTxtProgressBar(pb, 1)
return(results)
}

#####
### HAMILTONIAN MONTE CARLO - a step #####
#####

# How is it working? Take a look on a single step
HMC_step_norm <- function(epsilon, l, q.old, A){

  q = q.old
  p = rnorm(length(q), 0, 1) # independent standard normal variates
  p.old = p

  # Make half step for momentum at the beginning
  p = p - (epsilon / 2) * solve(A) %*% q

  # Alternate full steps for position and momentum
  for(i in 1:l){
    q.old <- q
    q = q + epsilon * p
    if(i == l) p = p - epsilon * solve(A) %*% q
    p.old = p

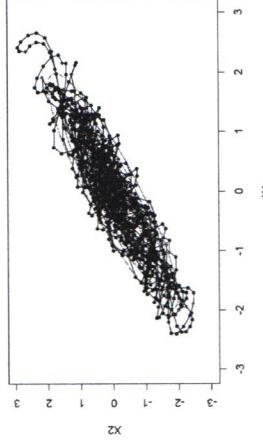
    # Evaluate potential and kinetic energies at start and end of trajectory
    current_U = t(q.old) %*% solve(A) %*% q.old / 2
    current_K = sum(p.old^2) / 2
    proposed_U = t(q) %*% solve(A) %*% q / 2
    proposed_K = sum(p^2) / 2

    # Evaluate potential and kinetic energies at start and end of trajectory
    current_U = t(q.old) %*% solve(A) %*% q.old / 2
    current_K = sum(p.old^2) / 2
    proposed_U = t(q) %*% solve(A) %*% q / 2
    proposed_K = sum(p^2) / 2

    # Accept or reject the state at end of trajectory, returning either
    # the position at the end of the trajectory or the initial position
    if(runif(1) < exp(current_U - proposed_U + current_K - proposed_K)){
      return(q) # accept
    } else{
      return(q.old) # reject
    }
  }
}

#####
# plot some iterations
q.old <- c(-1.5,-1)
A <- matrix(c(1,0,0,0,1,0), ncol = 2)

```



```
#####
##### HAMILTONIAN MONTE CARLO #####
#####

HMC_norm <- function(niter, burnin, x0, mu, Sigma, epsilon, L){
  results <- matrix(ncol = 2, nrow = niter - burnin)
  x.old <- x0

  pb <- txtProgressBar(0, niter, style = 3)

  for(i in 1:niter){

    x = x.old
    p = rnorm(length(x), 0, 1)
    P.old = p

    p = p - (epsilon / 2) * solve(Sigma) %% (x - mu)

    for (i in 1:L){
      x = x + epsilon * p
      if(i == L) p = p - epsilon * solve(Sigma) %% (x - mu)

      p = p - (epsilon / 2) * solve(Sigma) %% (x - mu)

      current_U = t(x.old - mu) %% solve(Sigma) %% (x.old - mu) / 2
      current_K = sum(P.old^2) / 2
      proposed_U = t(x - mu) %% solve(Sigma) %% (x - mu) / 2
      proposed_K = sum(t^2) / 2

      if (runif(1) < exp(current_U - proposed_U - current_K - proposed_K)){
        results[i - burnin, ] <- c(x.old)
      }
    }

    if(i > burnin){
      results[i - burnin, ] <- c(x.old)
    }

    if(runif(1) < exp(current_U - proposed_U - current_K - proposed_K)){
      x.old <- x
    }
  }

  setTxtProgressBar(pb, 1)
  close(pb)
  return(results)
}

# RW: HMC vs RWMH

```

```
HMC_sample <- HMC_norm(niter = niter, burnin = burnin, x0 = x0,
mu = mu, Sigma = Sigma, epsilon = 0.2, L = 10)
```

```
set.seed(123)
HMC_sample <- HMC_norm(niter = niter, burnin = burnin, x0 = x0,
mu = mu, Sigma = Sigma, epsilon = 0.2, L = 10)
```

```
HMC_sample <- RW_MH_bivariate(niter = niter, burnin = burnin, x0 = x0,
mu = mu, Sigma = Sigma)
```

Metropolis-Hastings
random-walk
(where the sampler is a bit slower
due to the ridge covariation)

```
#####
##### Diagnostic in CODA #####
#####

# import the results in coda (all together)
coda_ALL <- mcmc(cbind(HMC_sample, MH_sample), start = burnin + 1, end = niter)

# plot of the chains and summaries
plot(coda_ALL)
```

```
#####
##### plots #####
#####

# a common dataset for the two methods
pit_data <- data.frame(x1 = HMC_sample[,1], x2 = HMC_sample[,2],
y1 = MH_sample[,1], y2 = MH_sample[,2],
ind = 1:(niter - burnin))

# construct the data.frame for the contour plot
eg <- expand_grid(seq(-2, 6, length.out = 30), seq(-2, 6, length.out = 30))
dns_data <- data.frame(x = eg[,1], y = eg[,2],
z = apply(eg, 1, function(x) dmvnorm(x = x, mean = c(2, 2), sigma = A)))

pit_data <- pit_data[1:500,]

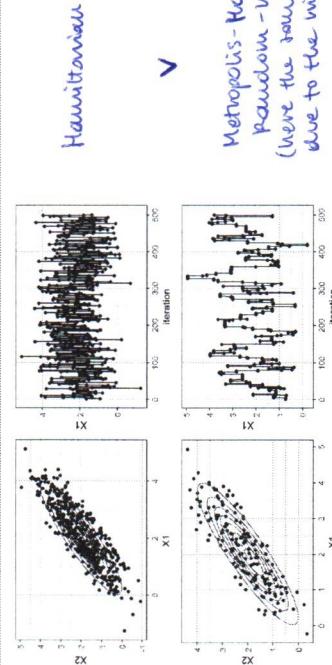
# plots for HMC
p1 <- ggplot(pit_data) +
  geom_point(aes(x = x1, y = x2)) +
  geom_contour(data = dns_data, aes(x = x, y = y, z = z)) +
  theme_bw() +
  xlab("X1") +
  ylab("X2")

p2 <- ggplot(pit_data) +
  geom_point(aes(x = x1, y = x2)) +
  geom_line(data = dns_data, aes(x = x, y = y, z = z)) +
  theme_bw() +
  xlab("X1") +
  ylab("X2")

# plots for MH
p3 <- ggplot(pit_data) +
  geom_point(aes(x = ind, y = y1)) +
  geom_contour(data = dns_data, aes(x = x, y = y, z = z)) +
  theme_bw() +
  xlab("X1") +
  ylab("X2")

p4 <- ggplot(pit_data) +
  geom_point(aes(x = ind, y = y1)) +
  geom_line(data = dns_data, aes(x = ind, y = y1)) +
  theme_bw() +
  xlab("Iteration") +
  ylab("X1") +
  xlab("X2") +
  ylab("X2")

# all the plots
grid.arrange(p1, p2, p3, p4, ncol = 2, nrow = 2)
```



```
#####
##### plots #####
#####

# plots
# a common dataset for the two methods
pit_data <- data.frame(x1 = HMC_sample[,1], x2 = HMC_sample[,2],
y1 = MH_sample[,1], y2 = MH_sample[,2],
ind = 1:(niter - burnin))

# construct the data.frame for the contour plot
eg <- expand_grid(seq(-2, 6, length.out = 30), seq(-2, 6, length.out = 30))
dns_data <- data.frame(x = eg[,1], y = eg[,2],
z = apply(eg, 1, function(x) dmvnorm(x = x, mean = c(2, 2), sigma = A)))

pit_data <- pit_data[1:500,]

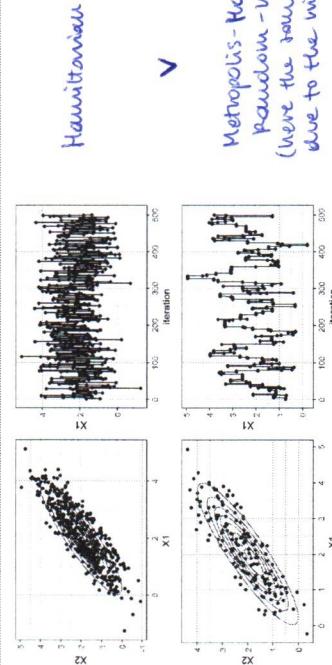
# plots for HMC
p1 <- ggplot(pit_data) +
  geom_point(aes(x = x1, y = x2)) +
  geom_contour(data = dns_data, aes(x = x, y = y, z = z)) +
  theme_bw() +
  xlab("X1") +
  ylab("X2")

p2 <- ggplot(pit_data) +
  geom_point(aes(x = x1, y = x2)) +
  geom_line(data = dns_data, aes(x = x, y = y, z = z)) +
  theme_bw() +
  xlab("X1") +
  ylab("X2")

# plots for MH
p3 <- ggplot(pit_data) +
  geom_point(aes(x = ind, y = y1)) +
  geom_contour(data = dns_data, aes(x = x, y = y, z = z)) +
  theme_bw() +
  xlab("X1") +
  ylab("X2")

p4 <- ggplot(pit_data) +
  geom_point(aes(x = ind, y = y1)) +
  geom_line(data = dns_data, aes(x = ind, y = y1)) +
  theme_bw() +
  xlab("Iteration") +
  ylab("X1") +
  xlab("X2") +
  ylab("X2")

# all the plots
grid.arrange(p1, p2, p3, p4, ncol = 2, nrow = 2)
```



```
#####
##### plots #####
#####

# plots
# a common dataset for the two methods
pit_data <- data.frame(x1 = HMC_sample[,1], x2 = HMC_sample[,2],
y1 = MH_sample[,1], y2 = MH_sample[,2],
ind = 1:(niter - burnin))

# construct the data.frame for the contour plot
eg <- expand_grid(seq(-2, 6, length.out = 30), seq(-2, 6, length.out = 30))
dns_data <- data.frame(x = eg[,1], y = eg[,2],
z = apply(eg, 1, function(x) dmvnorm(x = x, mean = c(2, 2), sigma = A)))

pit_data <- pit_data[1:500,]

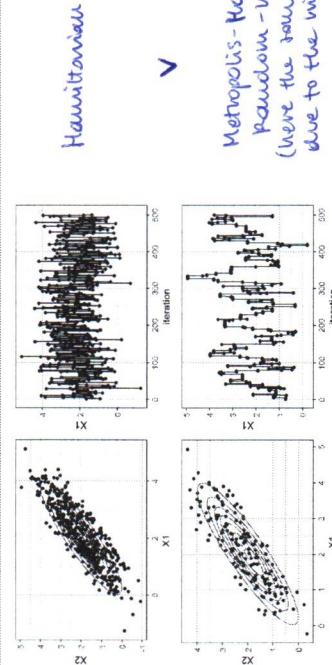
# plots for HMC
p1 <- ggplot(pit_data) +
  geom_point(aes(x = x1, y = x2)) +
  geom_contour(data = dns_data, aes(x = x, y = y, z = z)) +
  theme_bw() +
  xlab("X1") +
  ylab("X2")

p2 <- ggplot(pit_data) +
  geom_point(aes(x = x1, y = x2)) +
  geom_line(data = dns_data, aes(x = x, y = y, z = z)) +
  theme_bw() +
  xlab("X1") +
  ylab("X2")

# plots for MH
p3 <- ggplot(pit_data) +
  geom_point(aes(x = ind, y = y1)) +
  geom_contour(data = dns_data, aes(x = x, y = y, z = z)) +
  theme_bw() +
  xlab("X1") +
  ylab("X2")

p4 <- ggplot(pit_data) +
  geom_point(aes(x = ind, y = y1)) +
  geom_line(data = dns_data, aes(x = ind, y = y1)) +
  theme_bw() +
  xlab("Iteration") +
  ylab("X1") +
  xlab("X2") +
  ylab("X2")

# all the plots
grid.arrange(p1, p2, p3, p4, ncol = 2, nrow = 2)
```



```
#####
##### plots #####
#####

# plots
# a common dataset for the two methods
pit_data <- data.frame(x1 = HMC_sample[,1], x2 = HMC_sample[,2],
y1 = MH_sample[,1], y2 = MH_sample[,2],
ind = 1:(niter - burnin))

# construct the data.frame for the contour plot
eg <- expand_grid(seq(-2, 6, length.out = 30), seq(-2, 6, length.out = 30))
dns_data <- data.frame(x = eg[,1], y = eg[,2],
z = apply(eg, 1, function(x) dmvnorm(x = x, mean = c(2, 2), sigma = A)))

pit_data <- pit_data[1:500,]

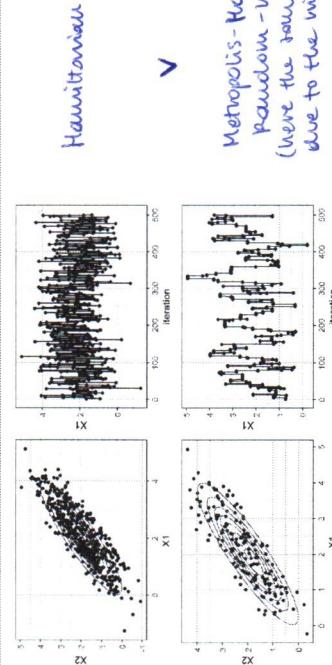
# plots for HMC
p1 <- ggplot(pit_data) +
  geom_point(aes(x = x1, y = x2)) +
  geom_contour(data = dns_data, aes(x = x, y = y, z = z)) +
  theme_bw() +
  xlab("X1") +
  ylab("X2")

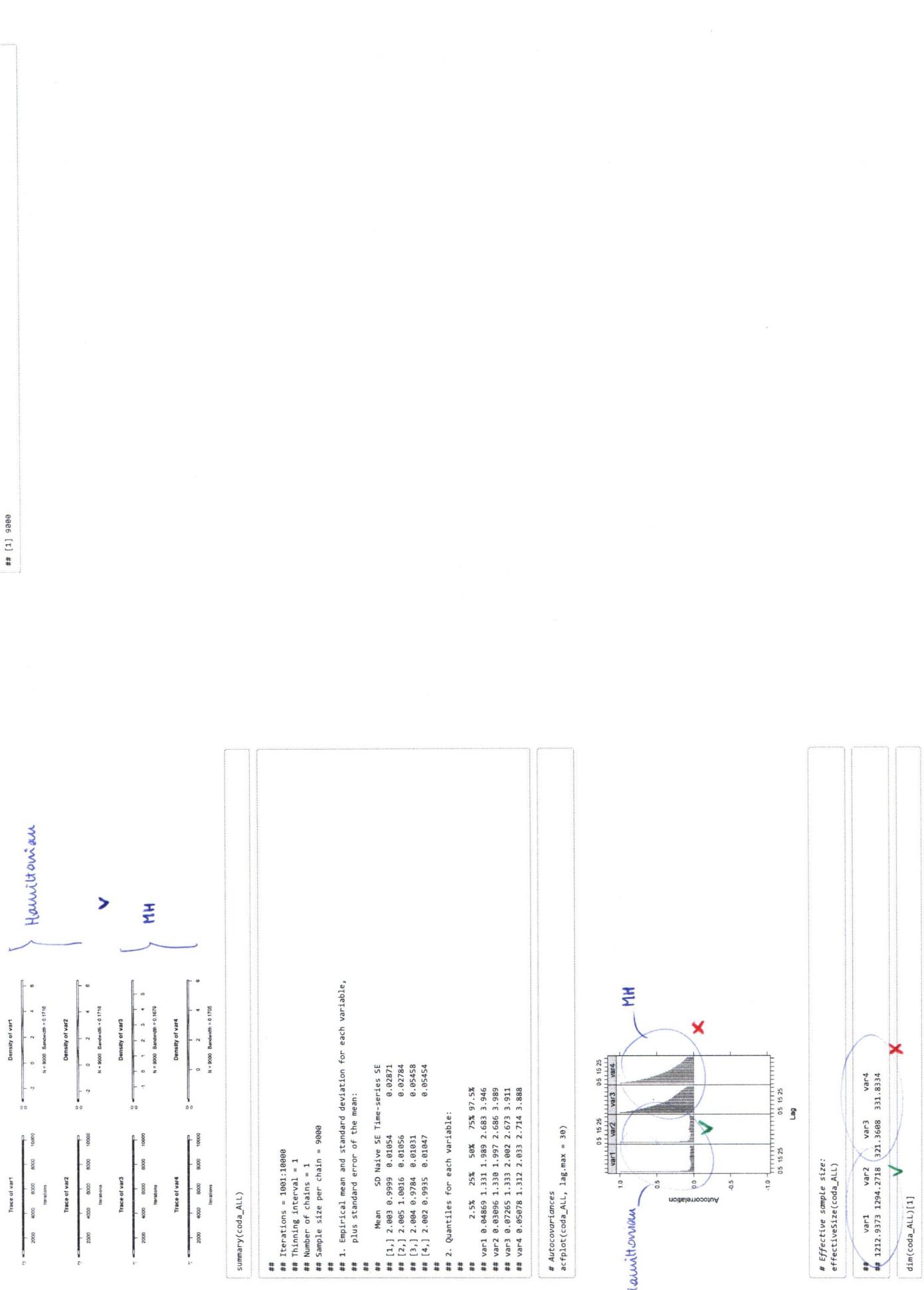
p2 <- ggplot(pit_data) +
  geom_point(aes(x = x1, y = x2)) +
  geom_line(data = dns_data, aes(x = x, y = y, z = z)) +
  theme_bw() +
  xlab("X1") +
  ylab("X2")

# plots for MH
p3 <- ggplot(pit_data) +
  geom_point(aes(x = ind, y = y1)) +
  geom_contour(data = dns_data, aes(x = x, y = y, z = z)) +
  theme_bw() +
  xlab("X1") +
  ylab("X2")

p4 <- ggplot(pit_data) +
  geom_point(aes(x = ind, y = y1)) +
  geom_line(data = dns_data, aes(x = ind, y = y1)) +
  theme_bw() +
  xlab("Iteration") +
  ylab("X1") +
  xlab("X2") +
  ylab("X2")

# all the plots
grid.arrange(p1, p2, p3, p4, ncol = 2, nrow = 2)
```





Introduction to STAN: hierarchical models and regression.

November 4, 2020

Gelman and Rubin Multiple Sequence Diagnostic

Idea: starting from several chains with different starting points, if the convergence is reached the trajectories should be similar.

Steps (for each parameter):

- 1) Run $m \geq 2$ chains of length $2n$ from overdispersed starting values
- 2) Discard the first n draws in each chain
- 3) Calculate the within-chain and between-chain variance
- 4) Calculate the estimated variance of the parameter as a weighted sum of the within-chain and between-chain variance
- 5) Calculate the potential scale reduction factor

in Stan is a popular thing to run more than one chain (and so we can do more diagnostic)

Note: here W , the within chain variance, is the mean of the variances of each chain and B , the between chain variance, is the variance of the chain means multiplied by n .

Gelman and Rubin Multiple Sequence Diagnostic²

We can estimate the variance of the stationary distribution as a weighted average of W and B .

The potential scale reduction factor is

$$\hat{R} = \sqrt{\frac{\text{Var}(\theta)}{W}}$$

when is ~ 1 is good because it means that all the chains look alike (all the chains look from the same distribution)

When \hat{R} is high (greater than 1.2) then we should run our chains out longer to improve convergence to the stationary distribution.

We can see how the potential scale reduction factor changes through the iterations using the `gelman.plot()` function.

A brief intro to STAN

<http://mc-stan.org/>

- STAN is a (very new!) probabilistic programming language implementing full Bayesian statistical inference with MCMC sampling (NUTS, HMC) and penalized maximum likelihood estimation with Optimization (BFGS).
- STAN is coded in C++ and runs on all major platforms (Linux, Mac, Windows).
- STAN is an open-source software developed mainly at Columbia University.
- STAN can be accessed through several interfaces: RStan (R), PyStan (Python), MatlabStan (Matlab) and also CmdStan (shell), <http://mc-stan.org/users/interfaces/cmdstan>

Sample from the posterior of the model

Multiple Markov chain Monte Carlo (MCMC) algorithms and optimization algorithms are implemented for the inference:

1 MCMC algorithms:

- Hamiltonian Monte Carlo (HMC) (default)
- No-U-Turn sampler

2 Optimization algorithms:

- BFGS algorithm (default), Nesterov's accelerated gradient descent algorithm, Newton's method

Typical workflow of using RStan

done with a simple command in R

- Represent a **statistical model** by writing its log-posterior density (up to a normalization constant independent from the parameters); this can be done using STAN modeling language
- Translate the model coded in **STAN** to C++ code using the function `stan`
- Compile the C++ code for the model using a C++ compiler (such as g++) to create a Dynamic Shared Object that can be loaded by R
- Run the DSO to **sample** from the posterior
- **Diagnose convergence** of the MCMC chains of sample
- **Conduct model inference**

we can just write the likelihood and the prior, stan will do the rest

[Don't worry! A single `rstan` call performs implicitly steps 2, 3 and 4.]

STAN has specific blocks...

- **DATA:**
 - Given input data
 - Executed first and load
- **TRANSFORMED DATA:**
 - Transform variables for convenience
- **PARAMETERS:**
 - Result output parameters
 - Updated at each iteration
- **TRANSFORM. PARAMETERS:**
 - Transform parameters for convenience
- **MODEL:**
 - Describe the model
- **GENERATED QUANTITIES:**
 - Generate quantities for monitoring convergence

... the order must be kept,
the blocks are optional (except model block)

describes the data, what we feed the model with

quantities of interest

} we specify the prior distribution and the likelihood

Variable and expression types

- **Primitive:** int and real
- **Matrix:** `vector[n]`, `row_vector[n]`, `matrix[m, n]`
- **Bounded:** primitive or matrix type, with
 $<lower = L>$, $<upper = U>$, $<lower = L, upper = U>$
- **Constrained:** `unit_vector` for unit-length vectors, `simplex` for unit simplexes, `ordered` for ordered vectors, `corr_matrix` and `cov_matrix` for symmetric and positive definite matrices.
- **Arrays:** collection of object of any type

Variable and expression types

- Sampling: $y \sim normal(mu, sigma)$
- Increments log-probability: `increment_log_prob(lp)`
add the value lp to the log-density
- For/while loop: `for(n in 1 : N), while(cond)`
- Block:{.....} (allows local variables)
- Print: `print(" TH = ", theta)`

Math functions

- All the built-in C + + functions and operators are available.
- For the matrices functions of basic arithmetics, solvers, decompositions, ... are available (check the manual).
- Each distribution of the library has:
 - pseudo random number generator
 - log density or mass function
 - cumulative distribution

Methods of the class S4 stanfit

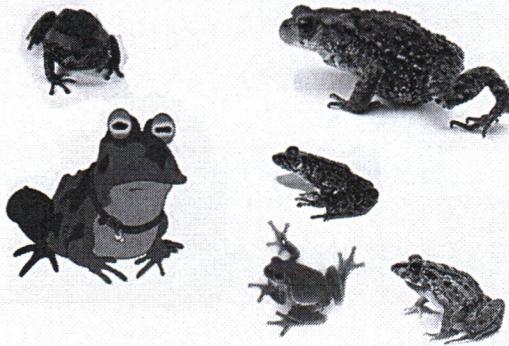
Printing, plotting, and summarizing:

- `show` Print the default summary for the model.
 - `print` Print a customizable summary for the model. See `print.stanfit`.
 - `plot` Create various plots summarizing the fitted model. See `plot.stanfit-method`.
 - `summary` Summarize the distributions of estimated parameters and derived quantities using the posterior draws. See `summary.stanfit-method`.
 - `get_posterior_mean` Get the posterior mean for parameters of interest (using `pars` to specify a subset of parameters). Returned is a matrix with one column per chain and an additional column for all chains combined.
-

see <https://mc-stan.org/rstan/reference/stanfit-class.html> for a summary

Example: hierarchical model

The data



politecnicimilano.webs.com is sharing your screen. Stop sharing Hide

12/19

The data²

We have a set of 341 observations divided into 12 genera. For each observation we have the body weight BoW, the length from the tail to the nose SVL and the brain weight BrW.

We are interested in studying a feature variable R which is measuring the relation between the body weight and the brain weight:

$$R = \log \left(\frac{BrW/BoW}{1 - (BrW/BoW)} \right)$$

Remark: $BrW \in (0, BoW)$, $BoW \in \mathbb{R}^+$, $R \in \mathbb{R}$

politecnicimilano.webs.com is sharing your screen. Stop sharing Hide

13/19

The model

Consider a hierarchical model as follows:

genus. stan

$$\begin{aligned} R_j &= R_1, \dots, R_{n_j} \mid \theta_j, \sigma^2 \sim N(\theta_j, \sigma^2), \quad j = 1, \dots, M \\ &\quad \text{group} \\ \theta_1, \dots, \theta_M &\mid \mu, \tau^2 \sim N(\mu, \tau^2) \\ (\mu, \tau^2) &\sim N(m_0, s_0^2) \times IG(\alpha_0, \beta_0) = \text{normal}(m_0, s_0^2) \cdot \text{inverse-gamma}(\alpha_0, \beta_0) \\ \sigma^2 &\sim IG(a_0, b_0) \end{aligned}$$

✓
II

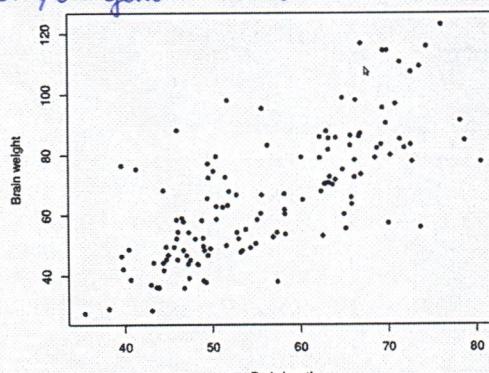
where the R_j 's are the brain/body weight indexes in the j -th genus, $\theta_1, \dots, \theta_M$ are the mean values for the index in the M groups, then we have the priors for the parameter (μ, τ^2) and for σ^2 .

The target distribution $X(\theta_1, \dots, \theta_M, \sigma^2, \mu, \tau^2 | R)$

14/19

The data

We want to study the relation Body-Brain weight in only one genus: "Rana"



politecnicimilano.webs.com is sharing your screen. Stop sharing Hide

15/19

The data²

We have a set of 133 observations for the *Rana* genus. For each observation we have the length from the tail to the nose X and the brain weight Y .

We are interested in studying the relation between the brain weight Y and the length of the body X , model of the kind $Y = g(X)$.

We consider two models:

- $Y = \beta_0 + \beta_1 X$
- $Y = \beta_0 + \beta_1 X + \beta_2 X^2$

16/19

The regression model I

Consider a **regression model** as follows:

$$\mathbf{Y} = Y_1, \dots, Y_n \mid \beta_0, \beta_1, \sigma^2 \sim N(\beta_0 + \beta_1 X, \sigma^2)$$

$$\beta_j \mid \tau_j^2 \sim N(0, \tau_j^2), \quad j = 1, 2$$

$$\sigma^2 \sim IG(a_0, b_0)$$

regression. stan



where the \mathbf{Y} are the brain weights, β_0, β_1 are the regression coefficients, then we have the prior for σ^2 .

The target distribution is $\mathcal{L}(\beta_0, \beta_1, \sigma^2, \mid \mathbf{Y}, \mathbf{X})$.

17/19

The regression model II

Consider a **regression model** as follows:

$$\mathbf{Y} = Y_1, \dots, Y_n \mid \beta_0, \beta_1, \beta_2, \sigma^2 \sim N(\beta_0 + \beta_1 X + \beta_2 X^2, \sigma^2)$$

$$\beta_j \mid \tau_j^2 \sim N(0, \tau_j^2), \quad j = 1, 2, 3$$

$$\sigma^2 \sim IG(a_0, b_0)$$

regression2. stan



where the \mathbf{Y} are the brain weights, $\beta_0, \beta_1, \beta_2$ are the regression coefficients, then we have the prior for σ^2 .

The target distribution is $\mathcal{L}(\beta_0, \beta_1, \beta_2, \sigma^2, \mid \mathbf{Y}, \mathbf{X})$.

18/19

WAIC

for
MODEL
COMPARISON

$$WAIC = -2 \text{ lppd} + 2p_{WAIC}$$

where lppd is the log pointwise predictive density

$$\begin{aligned} lppd &= \sum_{i=1}^n \log \left(\int p(Y_i \mid \theta) \pi(\theta \mid \mathbf{Y}) d\theta \right) \\ &\approx \sum_{i=1}^n \log \left(\frac{1}{S} \sum_{j=1}^S p(Y_i \mid \theta^{(j)}) \right) \end{aligned}$$

Log pointwise predictive density
(the integral is computed in
a Monte Carlo strategy)

and

$$\begin{aligned} p_{WAIC} &= \sum_{i=1}^n \text{Var}(\log p(Y_i \mid \theta) \mid \mathbf{Y}) \\ &\approx \sum_{i=1}^n \frac{1}{S-1} \sum_{j=1}^S \left(\log p(Y_i \mid \theta^{(j)}) - \overline{\log p(Y_i \mid \theta)} \right)^2 \end{aligned}$$

19/19

5/11/2020

4/11

```

#####
## regression.stan
#####

# we're giving to the model
# data
# {
#   int<lower = 0> M; // number of genus
#   int<lower = 0> N[M]; // numerosities in each genus: array of M cells
#   vector<lower = -25, upper = 25> sum(N); Y; // Logratio / (1 - ratio)
#   int<lower = 1, upper = M> genus_idx(sum(N)); // which genus the datum belongs to!
# }

# parameters
# {
#   real theta[M]; // array of the mean into each genus
#   real<lower = 0> sigma2; // inv_gamma(2., 1.) prior
#   real mu; // mean
#   real tau2; // std dev
# }

# model
# {
#   // Prior:
#   sigma2 ~ inv_gamma(2., 1.);
#   mu ~ normal(-5., 10.); // mean and std dev
#   tau2 ~ inv_gamma(2., 1.); // iid sample
#   theta ~ normal(mu, pow(tau2, 0.5));
#   // Likelihood:
#   for(i in 1:(sum(N)))
#   {
#     Y[i] ~ normal(theta[genus_idx[i]], pow(sigma2, 0.5));
#   }
# }

# generated quantities
# {
#   vector[N] log_lik;
#   for (j in 1:N)
#   {
#     log_lik[j] = normal_lpdf(Y[j] | mu[j], pow(sigma2, 0.5));
#   }
# }

```

```

#####
## regression2.stan
#####

# data
# {
#   # int<lower = 0> N;           // number of obs
#   # int<lower = 0> K;           // number of covariates (including the intercept)
#   # vector[N] Y;                // response
#   # matrix[N, K] X;             // covariates
#   #
#   // we pass also the parameters for the priors
#   # vector[K] mean_beta;
#   # vector[K] var_beta;
#   # real<lower = 0> scale_s2;
# }

# parameters
# {
#   # real<lower = 0> sigma2;
#   # vector[K] beta;
# }

# transformed parameters
# {
#   # vector[N] mu;
#   # mu = X * beta;
# }

# model
# {
#   // Prior:
#   # sigma2 ~ inv_gamma(2., scale_s2);
#   # for(k in 1:K)
#   {
#     # beta[k] ~ normal(mean_beta[k], pow(var_beta[k], 0.5));
#   }
#   #
#   // Likelihood:
#   # Y ~ normal(mu, pow(sigma2, 0.5));
# }

# generated quantities
# {
#   # vector[N] log_Lik;
#   # for (j in 1:N)
#   {
#     # log_Lik[j] = normal_lpdf(Y[j] | mu[j], pow(sigma2, 0.5));
#   }
# }

```

```
#####
##### EXAMPLE #####
##### Bayesian hierarchical model #####
#####

# R & STAN are friends!
library(rstan)

library(purrr)
library(coda)

# for plots
library(ggplot2)
library(tidyverse)

library(dplyr)

library(purrr)
library(coda)

# dataset: body weights and brain weights of frogs and salamander
# variable of interest  $\log(r / (1 - r))$  where  $r = \text{brain\_w} / \text{body\_w}$ 
Y <- read.csv("data.csv")
Y <- Y[complete.cases(Y)]
head(Y)

## # Order Family Genus Species Bod SVL BrW
## 1 Anura Bombinatoridae Bombina maxima 18543.9 60.24 45.30
## 2 Anura Bombinatoridae Bombina maxima 11781.3 61.54 52.98
## 3 Anura Bombinatoridae Bombina maxima 22459.8 59.80 58.70
## 4 Anura Bombinatoridae Bombina maxima 23549.8 60.39 33.80
## 5 Anura Bombinatoridae Bombina maxima 23891.0 60.51 40.00
## 6 Anura Bombinatoridae Bombina maxima 28676.2 61.89 41.00

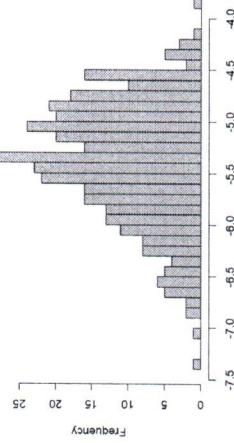
dim(Y)
## [1] 341 7

# Compute M, number of Genus, and the vector N,
# with the numerosities in each Genus:
m <- length(unique(Y[,3]))
num_Benus <- as.vector(table(Y[,3]))
table(Y[,3])

## # Amolops Bombina Bufo Chaperana Fejervarya Hyla
## 26 12 45 6 25 23
## Kaloula Narorana Paa Polypedates Rana Rhacophorus
## 17 6 6 16 133 26

genus_data <- list(M = m, N = num_Benus,
Y = log((Y[,7]/Y[,1]) / (1 - Y[,7])/Y[,5])),
genus_idx = as.numeric(as.factor(Y$genus)))
hist(genus_data$Y, breaks = 40, main = "", xlab = "")

genus_data <- list(M = m, N = num_Benus,
Y = log((Y[,7]/Y[,1]) / (1 - Y[,7])/Y[,5])),
genus_idx = as.numeric(as.factor(Y$genus)))
hist(genus_data$Y, breaks = 40, main = "", xlab = "")
```



```
#####
##### FIT THE MODEL #####
#####

# Fit a model defined in the stan modeling language
# and return the fitted result in an instance of class stanfit
# Steps performed by the function:
# (i) translate the model into C++ code
# (ii) compile the C++ code into a binary shared object
# (iii) sampling

?stan

## # starting httpd help server ...
## done

fit1 <- stan(file = "genus.stan",
data = genius.data,
iter = 5000, thin = 5,
chains = 2, warmup = 1000,
algorithm = 'NUTS', diagnostic_file = diagnostic.txt,
verbose = TRUE,
seed = 42) # To reproduce the results

## # Arguments:
## chains = number of chains that are randomly initialized, default is 4
## (iter, thin, warmup are setted for each chain)
## algorithm = 'NUTS', 'HMC'.
## sample_file, diagnostic_file = write files with samples or diagnostic data
## save_dso = whether save the dso object or not
## control = set the parameters of the algorithm and the adaptation parameters
is(fit1)
## [1] "stanfit"

## # We can split the 3 steps into 3 different functions:
## (i) stan: translate stan model into C++ code
step1 <- stan(file = 'genus.stan', model.name = 'genus_ex')
names(step1)

## [1] "status"   "model_cppname" "cpp_code"    "model_name"
## [5] "model_code"

## (ii) stan_model: construct a stan model that can be used to draw samples from the model
# The C++ code is compiled into a Dynamic Shared Object.
# DSO = A dynamic shared object (DSO) is an object file that is meant to be
# used simultaneously (or shared) by multiple applications while they are executing.
step2 <- stan_model(stan_ret = step1, verbose = FALSE)
is(step2)
```

```
# [1] "stanmodel"
```

```
# (i) sampling: draw samples from a stan-model
step3 <- sampling(step2,
  data = genus_data, chains = 2, iter = 2000,
  warmup = 500, thin = 1, seed = 42)
```

```
## SAMPLING FOR MODEL 'genus_ex' NOW (CHAIN 1).
## Chain 1:
## Gradient evaluation took 0 seconds
## 1000 transitions using 10 leapfrog steps per transition would take 0 seconds.
```

```
## Chain 1: Adjust your expectations accordingly!
## Chain 1:
## Chain 1: Iteration: 1 / 2000 [ 0% ] (warmup)
## Chain 1: Iteration: 200 / 2000 [ 10% ] (warmup)
## Chain 1: Iteration: 400 / 2000 [ 20% ] (warmup)
## Chain 1: Iteration: 501 / 2000 [ 25% ] (sampling)
## Chain 1: Iteration: 700 / 2000 [ 33% ] (sampling)
## Chain 1: Iteration: 900 / 2000 [ 43% ] (sampling)
## Chain 1: Iteration: 1100 / 2000 [ 53% ] (sampling)
## Chain 1: Iteration: 1300 / 2000 [ 63% ] (sampling)
## Chain 1: Iteration: 1500 / 2000 [ 73% ] (sampling)
## Chain 1: Iteration: 1700 / 2000 [ 83% ] (sampling)
## Chain 1: Iteration: 1900 / 2000 [ 93% ] (sampling)
## Chain 1: Iteration: 2000 / 2000 [100% ] (sampling)
## Chain 1:
## Elapsed Time: 0.229 seconds (warm-up)
## 0.469 seconds (sampling)
## 0.698 seconds (total)
```

```
## SAMPLING FOR MODEL 'genus_ex' NOW (CHAIN 2).
```

```
## Chain 2:
## Gradient evaluation took 0 seconds
## 1000 transitions using 10 leapfrog steps per transition would take 0 seconds.
```

```
## Chain 2: Adjust your expectations accordingly!
```

```
## Chain 2:
## Chain 2: Iteration: 1 / 2000 [ 0% ] (warmup)
## Chain 2: Iteration: 200 / 2000 [ 10% ] (warmup)
## Chain 2: Iteration: 400 / 2000 [ 20% ] (warmup)
## Chain 2: Iteration: 501 / 2000 [ 25% ] (sampling)
## Chain 2: Iteration: 700 / 2000 [ 33% ] (sampling)
## Chain 2: Iteration: 900 / 2000 [ 43% ] (sampling)
## Chain 2: Iteration: 1100 / 2000 [ 53% ] (sampling)
## Chain 2: Iteration: 1300 / 2000 [ 63% ] (sampling)
## Chain 2: Iteration: 1500 / 2000 [ 73% ] (sampling)
## Chain 2: Iteration: 1700 / 2000 [ 83% ] (sampling)
## Chain 2: Iteration: 1900 / 2000 [ 93% ] (sampling)
## Chain 2: Iteration: 2000 / 2000 [100% ] (sampling)
## Chain 2:
## Elapsed Time: 0.262 seconds (warm-up)
## 0.486 seconds (sampling)
## 0.748 seconds (total)
```

```
?sampling
```

```
is(step3)
```

```
# POST PROCESSING --
```

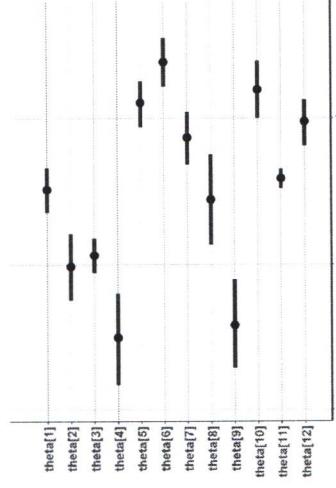
```
print(fit1,
probs = c(0.025, 0.5, 0.975),
par = c('mu', 'tau2', 'sigma2'))
```

```
## outer_level: 0.95 (95% intervals)
```

```
## c1_level: 0.95 (95% intervals)
```

```
# Draw the traceplot corresponding to one or more Markov chains,
# providing a visual way to inspect
# sampling behavior and assess mixing across chains and convergence.
```

```
rstan::traceplot(fit1, pars = c('mu', 'theta[1]'), inc_warmup = FALSE)
```



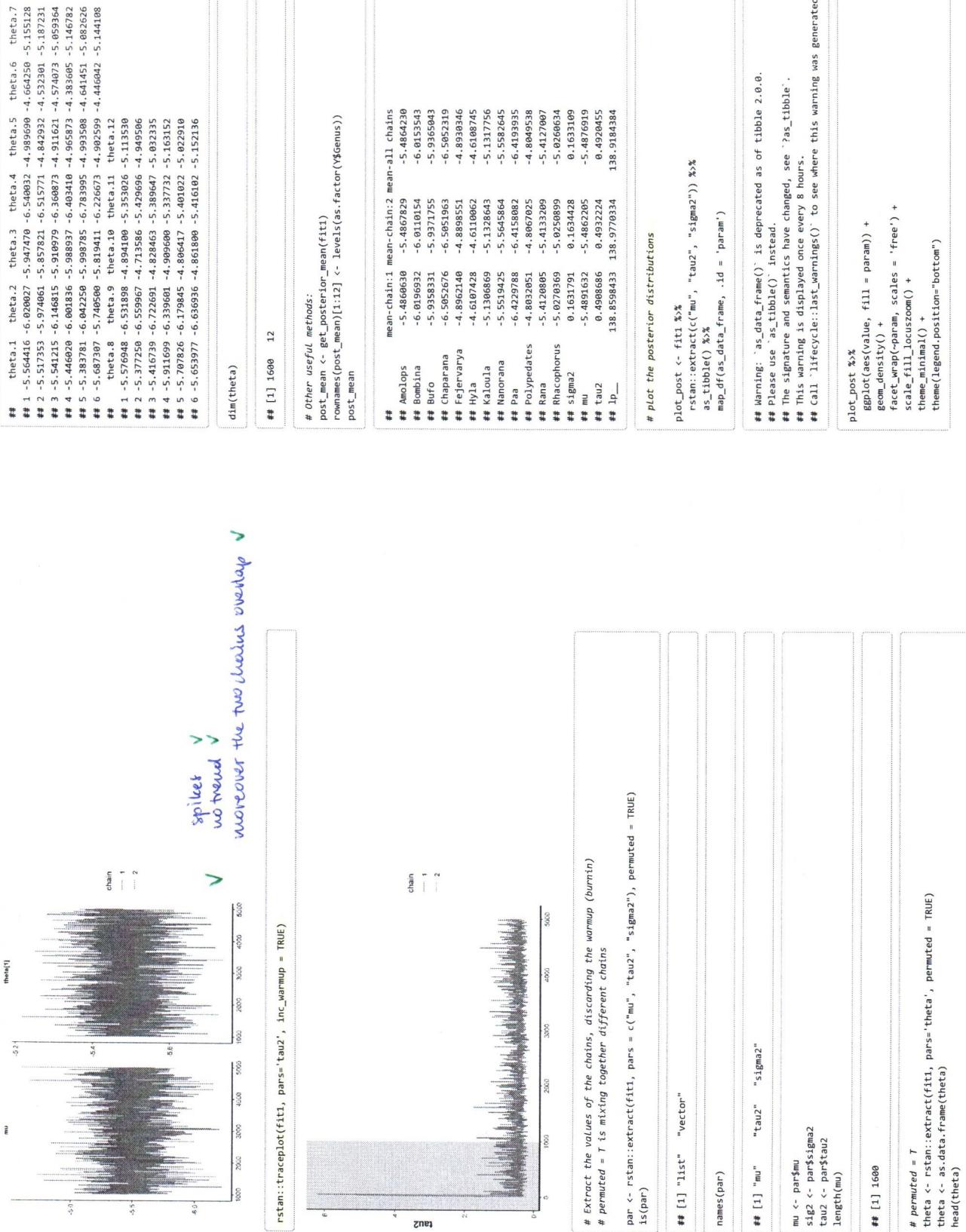
Output:

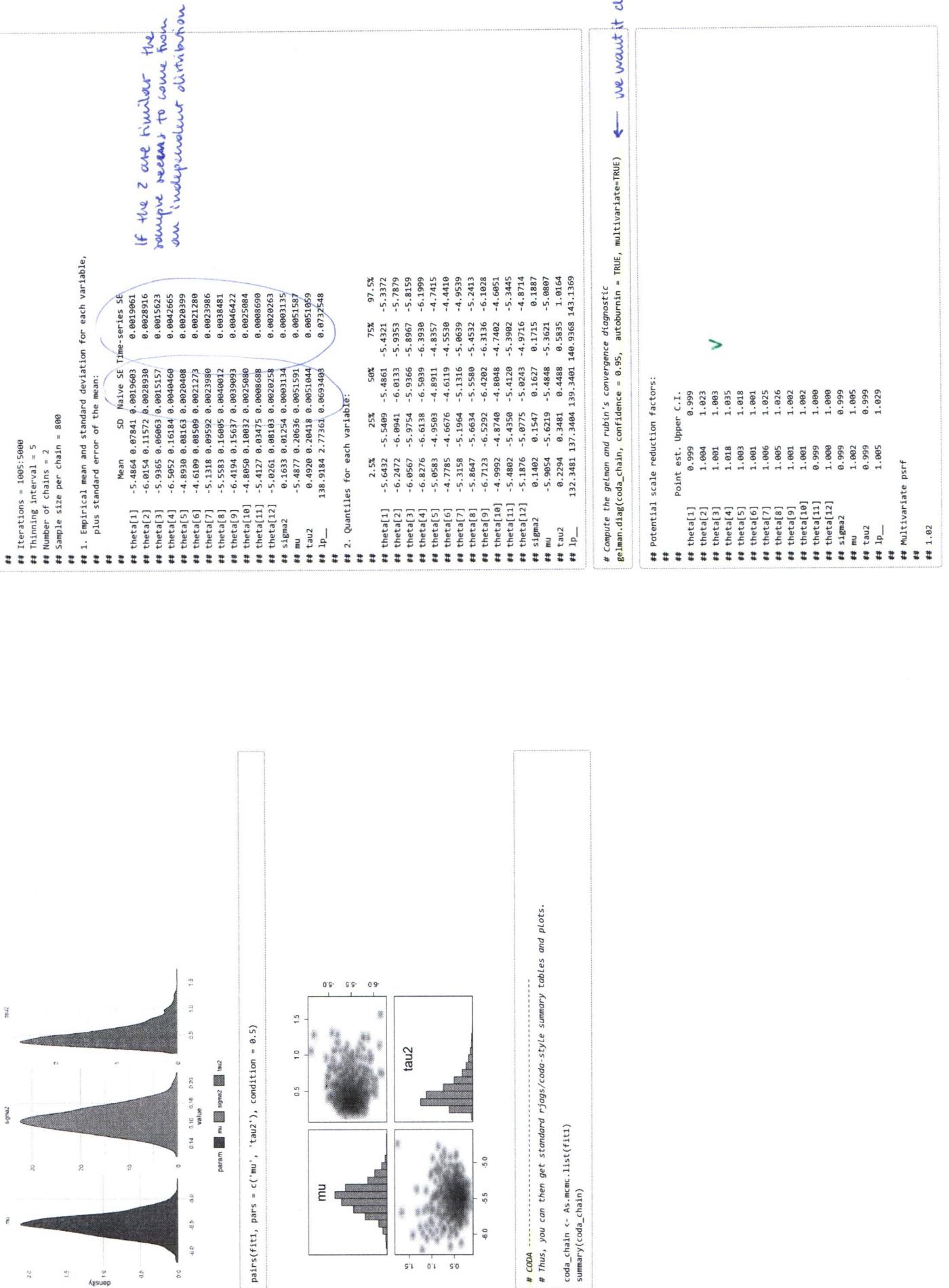
- Info about the iterations and the used sampler**
- Summary about parameters includes: the mean, the standard error of the mean ($\text{st.dev}/\sqrt{n}$), standard deviation and quantiles, the effective sample size and the split Rhat**
- All of them are computed without the warmup samples.**
- The split rate R:**
 - One way to monitor whether a chain has converged to the equilibrium distribution**
 - is to compare its behavior to other randomly initialized chains.**
 - R statistic measures the ratio of the average variance of samples within each chain to the variance of the pooled samples across chains; if all chains are at equilibrium, these will be the same and R will be one.**
 - If the chains have not converged to a common distribution, the R stat will be greater than 1.**

**mean (posterior) +
credible interval for
each parameter**

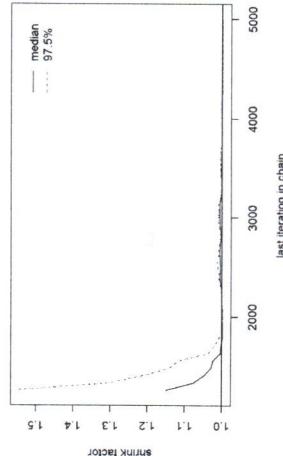
```
## Samples were drawn using NUTS(diag_e) at Sat Dec 26 00:30:08 2020.  
## For each parameter, n_eff is a crude measure of effective sample size.  
## and Rhat is the potential scale reduction factor on split chains (at  
## convergence, Rhat=1).
```

```
## Samples were drawn using NUTS(diag_e) at Sat Dec 26 00:30:08 2020.  
## For each parameter, n_eff is a crude measure of effective sample size.  
## and Rhat is the potential scale reduction factor on split chains (at  
## convergence, Rhat=1).
```





```
# & potential scale reduction factor is calculated for each variable in X,
# together with upper and lower confidence limits. Approximate convergence is
# diagnosed when the upper limit is close to 1. For multivariate chains, a multivariate
# value is calculated that bounds above the potential scale reduction factor for
# any linear combination of the (possibly transformed) variables
#
# This plot shows the evolution of Gelman and Rubin's shrink factor as the number
# of iterations increases
#
# gelman.diag(coda_chain, shrink=TRUE)
# gelman.plot(coda_chain[,1], confidence = 0.95)
```

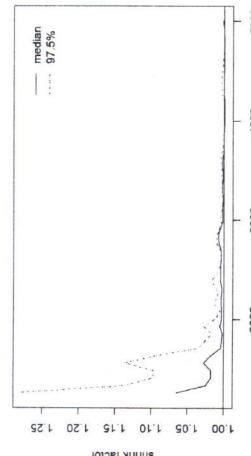


gelman.plot(coda_chain[,1], confidence = 0.95)

```
## [[1]]
## Fraction in 1st window = 0.1
## Fraction in 2nd window = 0.5
## theta[1] theta[2] theta[3] theta[4] theta[5] theta[6] theta[7] theta[8]
## 1.03137 -0.48831 -0.51213 2.00768 1.62820 -0.08911 0.06533 -0.31644
## theta[9] theta[10] theta[11] theta[12] sigma2 mu tau2 lp_
## -1.29011 2.55167 0.95158 0.19530 -1.19261 0.75953 0.37998 0.18273
## [[2]]
## Fraction in 1st window = 0.1
## Fraction in 2nd window = 0.5
## theta[1] theta[2] theta[3] theta[4] theta[5] theta[6] theta[7] theta[8]
## -2.23974 1.87965 0.62554 -0.53404 2.06976 2.06815 1.26767 0.39706
## theta[9] theta[10] theta[11] theta[12] sigma2 mu tau2 lp_
## 0.88466 0.52481 -1.98652 -0.45858 1.00894 -0.30376 0.48857 1.10584
## OPTIMIZATION: obtain a point estimate by maximizing the joint posterior
# from the model defined by class stanmodel
is(step2)
```

```
opt <- optimizing(step2,
  data = genus_data,
  algorithm = 'Newton')
opt
```

```
## $par
## theta[1] theta[2] theta[3] theta[4] theta[5] theta[6] theta[7]
## -5.4859567 -6.0162279 -5.9378624 -6.4961737 -4.8988234 -4.6476175 -5.138464
## theta[8] theta[9] theta[10] theta[11] theta[12] sigma2 mu
## -5.5575901 -6.4089653 -4.88952538 -5.4125007 -5.0287780 0.1559170 -5.4820056
## tau2
## 0.3478813
## $value
## [1] 149.4373
## $return_code
## [1] 0
## $theta_filde
## theta[1] theta[2] theta[3] theta[4] theta[5] theta[6] theta[7]
## [1,] -5.485957 -6.016228 -5.937862 -4.898823 -4.64762 -5.13846
## [2,] theta[8] theta[9] theta[10] theta[11] theta[12] sigma2 mu
## [3,] -5.55759 -6.408965 -4.8895254 -5.412501 -5.028778 0.155917 -5.482006
## tau2
## [1] 0.3478813
```



```
# Geweke's convergence diagnostic
Geweke.diag(coda_chain, frac1=0.1, frac2=0.5)
```

```

# Arguments:
#   algorithm = 'LBFGS', 'BFGS', 'Newton'

# RUN MULTIPLE CHAINS IN PARALLEL - -----
# Specifying the argument chains, rstan runs chains sequentially using one R
# process, which means that rstan does not support sampling in parallel
# directly. However, the function SFLIST2STANFIT merges a list of stanfit objects
# into one.

# Argument: a list of stanfit objects. The stanfit objects to be merged
# need to have the same configuration of iteration, warmup and thinning
# Output: an object of stanfit consolidated from all the input stanfit objects
# We can run multiple chains in parallel using using other approaches or packages
# and then use this function to merge them.

library(parallel)

# Create the model:
fit2 <- stan_model(file = 'genus.stan')

# Set parameters of the chain:
IT = 1000
WJ = 1000
TH = 5
seed <- 42

# Note: the following does not work on windows systems (use function parLapply or
# package foreach)
# mc.cores = The number of cores to use, i.e. at most how many child processes will
# be run simultaneously. The option is initialized from environment variable MC_CORES
# if set. Must be at least one, and parallelization requires at least two cores.
detectcores()

# # multicore
# system.time(
#   sflist <- mclapply(1:4,
#                     function(i)
#                       sampling(step2,
#                               data = genus_data, chains = 1, iter = IT*TH+i*WJ,
#                               warmup = WJ, thin=TH,
#                               chain_id = i,
#                               refresh = -1,
#                               pars = c('mu', 'sigma2', 'tau2'))),
#   mc.cores = 4
# )

# single core
system.time(
  sflist <- mclapply(1:4,
                    function(i)
                      sampling(step2,
                               data = genus_data, chains = 1, iter = IT*TH+i*WJ,
                               warmup = WJ, thin=TH,
                               chain_id = i,
                               refresh = -1,
                               pars = c('mu', 'sigma2', 'tau2'))),
  mc.cores = 1
)

## Chain 1:
## Chain 1: Gradient evaluation took 0 seconds
## Chain 1: 1000 transitions using 10 leapfrog steps per transition would take 0 seconds.
## Chain 1: Adjust your expectations accordingly!
## Chain 1: Elapsed Time: 0.468 seconds (Sampling)
## Chain 1: 1.883 seconds (Warm-up)
## Chain 1: 2.351 seconds (Total)

## Chain 2:
## Chain 2: Gradient evaluation took 0 seconds
## Chain 2: 1000 transitions using 10 leapfrog steps per transition would take 0 seconds.
## Chain 2: Adjust your expectations accordingly!
## Chain 2: Elapsed Time: 0.64 seconds (Sampling)
## Chain 2: 1.924 seconds (Warm-up)
## Chain 2: 2.564 seconds (Total)

## Chain 3:
## Chain 3: Gradient evaluation took 0 seconds
## Chain 3: 1000 transitions using 10 leapfrog steps per transition would take 0 seconds.
## Chain 3: Adjust your expectations accordingly!
## Chain 3: Elapsed Time: 0.414 seconds (Sampling)
## Chain 3: 1.536 seconds (Warm-up)
## Chain 3: 1.95 seconds (Total)

## Chain 4:
## Chain 4: Gradient evaluation took 0.001 seconds
## Chain 4: 1000 transitions using 10 leapfrog steps per transition would take 10 seconds.
## Chain 4: Adjust your expectations accordingly!
## Chain 4: Elapsed Time: 0.391 seconds (Sampling)
## Chain 4: 2.039 seconds (Warm-up)
## Chain 4: 2.43 seconds (Total)

is(sflist)
## [1] "list"  "vector"

## Summary for 4 different chains
print(sflist)

```

```

## [[1]]
## Inference for Stan model: genus_ex.
## 1 chains, each with iter=6000; warmup=1000; thin=5;
## post-warmup draws per chain=1000, total post-warmup draws=6000.
## 
## mean se_mean sd 2.5% 50% 75% 97.5% n_eff Rhat
## mu -5.49 0.01 0.22 -5.93 -5.63 -5.36 -5.02 839 1
## sigma2 0.16 0.00 0.01 0.14 0.15 0.16 0.17 0.19 1101 1
## tau2 0.59 0.01 0.23 0.24 0.36 0.45 0.60 1.05 993 1
## lp__ 138.91 0.09 2.81 132.41 137.39 139.22 140.98 143.41 999 1
## 
## Samples were drawn using NUTS(diag_e) at Sat Dec 26 00:31:13 2020.
## For each parameter, n_eff is a crude measure of effective sample size,
## and Rhat is the potential scale reduction factor on split chains (at
## convergence, Rhat=1).
## 

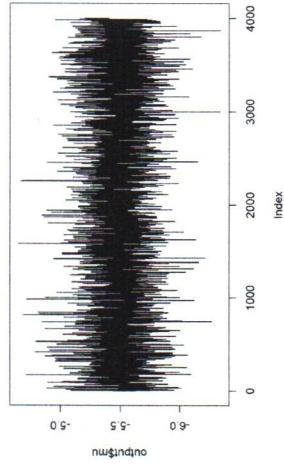
[[2]]
## Inference for Stan model: genus_ex.
## 1 chains, each with iter=6000; warmup=1000; thin=5;
## post-warmup draws per chain=1000, total post-warmup draws=6000.
## 
## mean se_mean sd 2.5% 50% 75% 97.5% n_eff Rhat
## mu -5.48 0.01 0.20 -5.67 -5.61 -5.49 -5.36 -5.09 938 1
## sigma2 0.16 0.00 0.01 0.14 0.15 0.16 0.17 0.19 1037 1
## tau2 0.49 0.01 0.21 0.22 0.35 0.45 0.58 1.04 957 1
## lp__ 139.04 0.09 2.72 133.08 137.21 139.46 141.09 143.46 870 1
## 
## Samples were drawn using NUTS(diag_e) at Sat Dec 26 00:31:16 2020.
## For each parameter, n_eff is a crude measure of effective sample size,
## and Rhat is the potential scale reduction factor on split chains (at
## convergence, Rhat=1).
## 

[[3]]
## Inference for Stan model: genus_ex.
## 1 chains, each with iter=6000; warmup=1000; thin=5;
## post-warmup draws per chain=1000, total post-warmup draws=6000.
## 
## mean se_mean sd 2.5% 50% 75% 97.5% n_eff Rhat
## mu -5.48 0.01 0.20 -5.66 -5.61 -5.47 -5.34 -5.10 759 1
## sigma2 0.16 0.00 0.01 0.14 0.15 0.16 0.17 0.19 882 1
## tau2 0.49 0.01 0.22 0.22 0.34 0.43 0.58 1.04 950 1
## lp__ 138.93 0.09 2.76 132.77 137.21 139.36 140.99 143.49 931 1
## 
## Samples were drawn using NUTS(diag_e) at Sat Dec 26 00:31:18 2020.
## For each parameter, n_eff is a crude measure of effective sample size,
## and Rhat is the potential scale reduction factor on split chains (at
## convergence, Rhat=1).
## 

[[4]]
## Inference for Stan model: genus_ex.
## 1 chains, each with iter=6000; warmup=1000, total post-warmup draws=6000.
## 
## mean se_mean sd 2.5% 50% 75% 97.5% n_eff Rhat
## mu -5.49 0.01 0.21 -5.89 -5.63 -5.48 -5.35 -5.07 1134 1
## sigma2 0.16 0.00 0.01 0.14 0.15 0.16 0.17 0.19 1053 1
## tau2 0.59 0.01 0.23 0.23 0.36 0.46 0.58 1.10 1109 1
## lp__ 138.96 0.09 2.79 132.79 137.26 139.27 140.89 143.41 907 1
## 
## Samples were drawn using NUTS(diag_e) at Sat Dec 26 00:31:20 2020.
## For each parameter, n_eff is a crude measure of effective sample size,
## and Rhat is the potential scale reduction factor on split chains (at
## convergence, Rhat=1).

# Merge the chains: we will obtain IT*NUM_CORES iterations
fit_p <- sflist2stanfit(sflist)
print(fit_p)

```



```

# plot the posterior distributions
plot_post <- output %>%
  as.data.frame() %>%
  map_df(as_data.frame, .id = 'param')

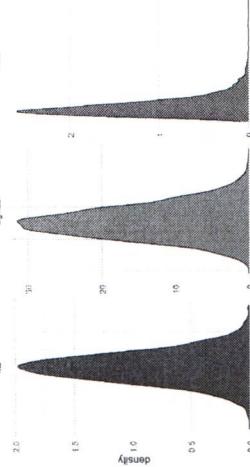
plot_post %>%
  ggplot(aes(value, fill = param)) +
  geom_density() +
  facet_wrap(~param, scales = 'free') +
  scale_fill_lowuscom() +
  theme_minimal() +
  theme(legend.position="bottom")

```

```

## SAMPLING FOR MODEL 'regression' NOW (CHAIN 1).
## Chain 1:
## Gradient evaluation took 0 seconds
## 1000 transitions using 10 leapfrog steps per transition would take 0 seconds.
## Chain 1: Adjust your expectations accordingly!
## Chain 1:
## Chain 1: Iteration: 1 / 5000 [ 0%] (warmup)
## Chain 1: Iteration: 500 / 5000 [ 10%] (warmup)
## Chain 1: Iteration: 1000 / 5000 [ 20%] (warmup)
## Chain 1: Iteration: 1001 / 5000 [ 20%] (sampling)
## Chain 1: Iteration: 1500 / 5000 [ 30%] (sampling)
## Chain 1: Iteration: 2000 / 5000 [ 40%] (sampling)
## Chain 1: Iteration: 2500 / 5000 [ 50%] (sampling)
## Chain 1: Iteration: 3000 / 5000 [ 60%] (sampling)
## Chain 1: Iteration: 3500 / 5000 [ 70%] (sampling)
## Chain 1: Iteration: 4000 / 5000 [ 80%] (sampling)
## Chain 1: Iteration: 4500 / 5000 [ 90%] (sampling)
## Chain 1: Iteration: 5000 / 5000 [100%] (sampling)
## Chain 1: Elapsed Time: 0.132 seconds (Warm-up)
## Chain 1: 0.456 seconds (Sampling)
## Chain 1: 0.588 seconds (Total)

## SAMPLING FOR MODEL 'regression' NOW (CHAIN 2).
## Chain 2:
## Gradient evaluation took 0.001 seconds
## 1000 transitions using 10 leapfrog steps per transition would take 10 seconds.
## Chain 2: Adjust your expectations accordingly!
## Chain 2:
## Chain 2: Iteration: 1 / 5000 [ 0%] (warmup)
## Chain 2: Iteration: 500 / 5000 [ 10%] (warmup)
## Chain 2: Iteration: 1000 / 5000 [ 20%] (sampling)
## Chain 2: Iteration: 1001 / 5000 [ 20%] (sampling)
## Chain 2: Iteration: 1500 / 5000 [ 30%] (sampling)
## Chain 2: Iteration: 2000 / 5000 [ 40%] (sampling)
## Chain 2: Iteration: 2500 / 5000 [ 50%] (sampling)
## Chain 2: Iteration: 3000 / 5000 [ 60%] (sampling)
## Chain 2: Iteration: 3500 / 5000 [ 70%] (sampling)
## Chain 2: Iteration: 4000 / 5000 [ 80%] (sampling)
## Chain 2: Iteration: 4500 / 5000 [ 90%] (sampling)
## Chain 2: Iteration: 5000 / 5000 [100%] (sampling)
## Chain 2: Elapsed Time: 0.133 seconds (Warm-up)
## Chain 2: 0.444 seconds (Sampling)
## Chain 2: 0.577 seconds (Total)
## Chain 2:
```

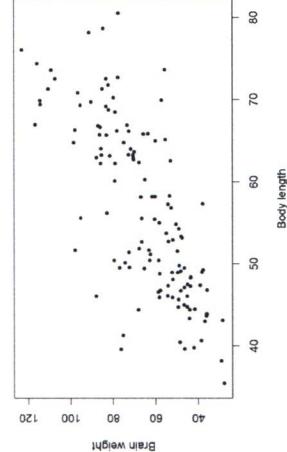


```

#####
##### EXAPLIE #####
##### Bayesian regression #####
#####

regression_data <- list(N = sum(Y$genus == "Rana"),
Y = Y$BML[Y$genus == "Rana"],
X = Y$VL[X$genus == "Rana"])
plot(regression_data$x, regression_data$y, pch = 20,
xlab = "Body length", ylab = "Brain weight")

```



```
# FIT THE MODEL
```

```

fit2 <- stanfile = "regression.stan",
      data = regression_data,
      iter = 5000, thin = 5,
      chains = 2, warmup = 1000,
      algorithm = 'NUTS',
      diagnostic_file = 'diag_reg.txt',
      verbose = FALSE,
      seed = 42)

```

```
print(fit2, par = c('beta0', 'beta1', 'sigma2'))
```

```

## Inference for Stan model: regression.
## 2 chains, each with iter=5000; warmup=1000; thin=5;
## post-warmup draws per chain=800, total post-warmup draws=1600.
## 
## mean se_mean sd 2.5% 50% 75% 97.5% n_eff Rhat
## beta0 -0.37 0.02 0.97 -1.02 -0.35 0.29 1.48 1683 1
## beta1 1.18 0.00 0.63 1.13 1.17 1.18 1.20 1.24 1577 1
## sigma2 269.07 0.67 25.55 164.93 191.20 267.29 224.52 264.24 1473 1

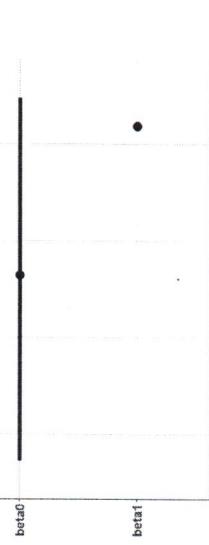
## Samples were drawn using NUTS(diag.e) at Sat Dec 26 00:32:14 2020.
## For each parameter, n_eff is a crude measure of effective sample size,
## and Rhat is the potential scale reduction factor on split chains (at
## convergence, Rhat=1).
```

```

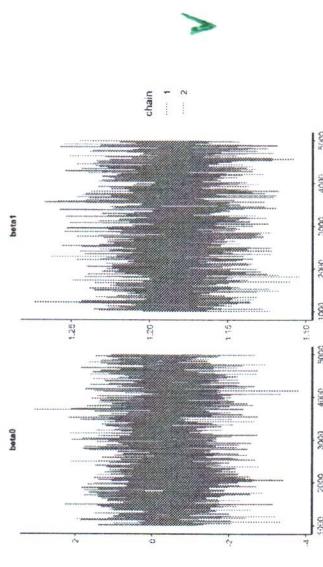
## ci_level: 0.95 (95% intervals)
## outer_level: 0.95 (95% intervals)

plot(fit2, ask = T, pars = c("beta0", "beta1", ci_level = 0.95, fill_color = "blue"))

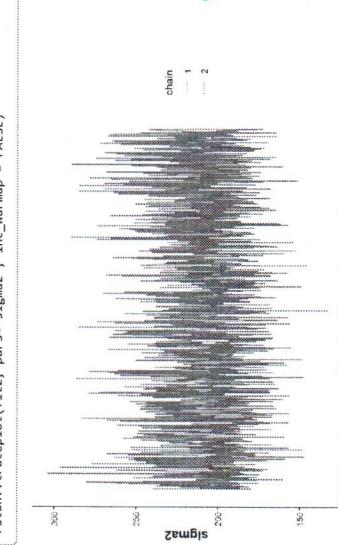
```



```
# Draw the traceplot corresponding to one or more Markov chains,
# providing a visual way to inspect
# sampling behavior and assess mixing across chains and convergence.
rstan::traceplot(fit2, pars = c("beta0", "beta1"), inc_warmup = FALSE)
```



```
rstan::traceplot(fit2, pars="sigma2", inc_warmup = FALSE)
```



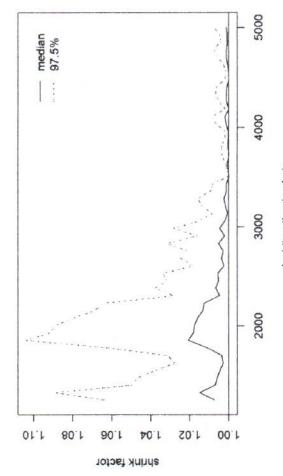
```
#####
# DIAGNOSTIC IN CODA #####
# --#
coda_chain_reg <- As.mcmc.list(fit2, pars = c("beta0", "beta1", "sigma2"))
summary(coda_chain_reg)
```

```
## Iterations = 1005:5000
## Thinning interval = 5
## Number of chains = 2
## Sample size per chain = 800
## 
## 1. Empirical mean and standard deviation for each variable,
## plus standard error of the mean:
## 
##           Mean        SD  Naive SE Time-series SE
## beta0    -0.3662  0.9727  0.0243167   0.0243234
## beta1    1.1841  0.0268  0.0006599   0.0007176
## sigma2 299.0738 25.5544  0.6388681   0.6388625
## 
## 2. Quantiles for each variable:
## 
##          2.5%     25%      50%     75%    97.5%
## beta0  -2.270  -1.819  -0.3539  0.2915  1.477
## beta1   1.132  1.165  1.1837  1.2026  1.237
## sigma2 164.927 191.196 267.2946 224.5227 264.237
## 
## Compute the gelman and rubin's convergence diagnostic
gelman.diag(coda_chain_reg, confidence = 0.95, autoburnin = TRUE, multivariate=TRUE)
```

```
## Potential scale reduction factors:
## 
##          Point est. Upper C.I.
## beta0      1       1.01
## beta1      1       1.00
## sigma2     1       1.00
## 
## Multivariate psrf
## 
##          1
```

```
# The potential scale reduction factor(22) is calculated for each variable in x,
# together with upper and lower confidence limits. Approximate convergence is
# diagnosed when the upper limit is close to 1. For multivariate chains, a multivariate
# value is calculated that bounds above the potential scale reduction factor for
# any linear combination of the (possibly transformed) variables
```

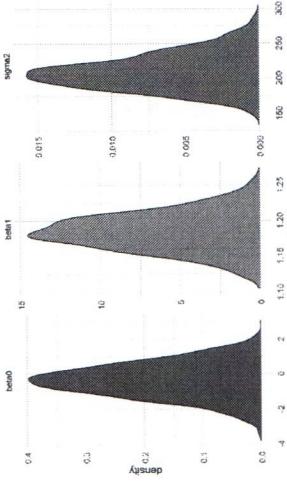
```
# This plot shows the evolution of Gelman and Rubin's shrink factor as the number
# of iterations increases
gelman.plot(coda_chain_reg[,1], confidence = 0.95)
```



```
gelman.plot(coda_chain_reg[,2], confidence = 0.95)
```

```
# plots of the posteriors
plot_post <- fit2 %>%
  rstan::extract(c("beta0", "beta1", "sigma2")) %>%
  as.data.frame() %>%
  map_df(as_data_frame, .id = 'param')

plot_post %>%
  ggplot(aes(value, fill = param)) +
  geom_density() +
  facet_wrap(~param, scales = 'free') +
  scale_fill_locuszoom() +
  theme_minimal() +
  theme(legend.position="bottom") +
  labs(fill = "")
```



```
#####
##### EXAMPLE #####
##### Bayesian regression (square) #####
#####

regression_data_full <- list(N = sum(YGenus == "Rana"),
                                K = 3,
                                Y = YBr-W[YGenus == "Rana"], 
                                X = cbind(rep(1, sum(YGenus == "Rana")),
                                          YSL[YGenus == "Rana"],
                                          (YSL[YGenus == "Rana"])^2),
                                scale_s2 = 10)

reg_lm <- summary(lm(regression_data_full$Y ~ regression_data_full[,2:3]))
regression_data_full$mean_beta <- reg_lm$coefficients[,1]
regression_data_full$par_beta <- (reg_lm$coefficients[,2])^2

# FIT THE MODEL (again) ----
fit3 <- stan(file = "regression2.stan",
             data = regression_data_full,
             iter = 5000, thin = 5,
             chains = 2, warmup = 1000,
             algorithm = 'NUTS',
             diagnostic_file = 'diag_reg.txt',
             verbose = FALSE,
             seed = 42)
```



```
# Gelman's convergence diagnostic
```

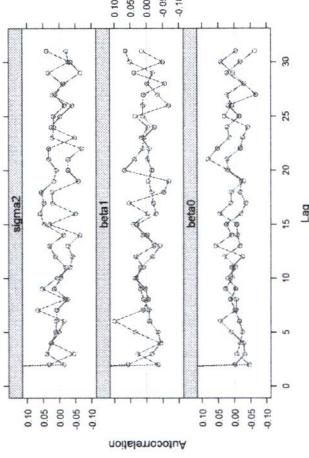
```
gelman.diag(coda_chain_Reg, frac1=0.1, frac2=0.5)
```

```
## [[1]]
## Fraction in 1st window = 0.1
## Fraction in 2nd window = 0.5
## beta0
## beta1 sigma2
## -1.632 1.707 4.038
## 
```

```
## [[2]]
## Fraction in 1st window = 0.1
## Fraction in 2nd window = 0.5
## beta0
## beta1 sigma2
## 0.5899 -0.4188 1.6639
```

```
# autocorrelation and plot
```

```
acfplot(coda_chain_Reg, lag.max = 30)
```



```

## # # SAMPLING FOR MODEL 'regression2' NOW (CHAIN 1).
## # # Chain 1:
## # # Chain 1: Gradient evaluation took 0 seconds
## # # Chain 1: 1000 steps using 10 leapfrog steps per transition would take 0 seconds.
## # # Chain 1: Adjust your expectations accordingly!
## # # Chain 1:
## # # Chain 1: Iteration: 1 / 5000 [ 0%] (warmup)
## # # Chain 1: Iteration: 500 / 5000 [ 10%] (warmup)
## # # Chain 1: Iteration: 1000 / 5000 [ 20%] (warmup)
## # # Chain 1: Iteration: 1001 / 5000 [ 20%] (sampling)
## # # Chain 1: Iteration: 1500 / 5000 [ 30%] (sampling)
## # # Chain 1: Iteration: 2000 / 5000 [ 40%] (sampling)
## # # Chain 1: Iteration: 2500 / 5000 [ 50%] (sampling)
## # # Chain 1: Iteration: 3000 / 5000 [ 60%] (sampling)
## # # Chain 1: Iteration: 3500 / 5000 [ 70%] (sampling)
## # # Chain 1: Iteration: 4000 / 5000 [ 80%] (sampling)
## # # Chain 1: Iteration: 4500 / 5000 [ 90%] (sampling)
## # # Chain 1: Iteration: 5000 / 5000 [100%] (sampling)
## # # Chain 1:
## # # Chain 1: Elapsed Time: 1.793 seconds (warm-up)
## # # Chain 1: 4.193 seconds (Sampling)
## # # Chain 1: 5.986 seconds (Total)

```

session? Now (Chain 2).

```

## SAMPLING FOR MODEL 'regression2' NOW ((CHAIN 2)).

## Chain 2:
## Chain 2: Gradient evaluation took 0 seconds
## Chain 2: 1000 transitions using 10 leapfrog steps per transition would take 0 seconds.
## Chain 2: Adjust your expectations accordingly!
## Chain 2:
## Chain 2: Iteration: 1 / 5000 [  0%] (warmup)
## Chain 2: Iteration: 500 / 5000 [ 10%] (warmup)
## Chain 2: Iteration: 1000 / 5000 [ 20%] (warmup)
## Chain 2: Iteration: 1500 / 5000 [ 30%] (Sampling)
## Chain 2: Iteration: 1980 / 5000 [ 39%] (Sampling)
## Chain 2: Iteration: 2000 / 5000 [ 40%] (Sampling)
## Chain 2: Iteration: 2500 / 5000 [ 50%] (Sampling)
## Chain 2: Iteration: 3000 / 5000 [ 60%] (Sampling)
## Chain 2: Iteration: 35000 / 5000 [ 70%] (Sampling)
## Chain 2: Iteration: 40000 / 5000 [ 80%] (Sampling)
## Chain 2: Iteration: 45000 / 5000 [ 90%] (Sampling)
## Chain 2: Iteration: 50000 / 5000 [100%] (Sampling)

## Chain 2:
## Chain 2: Elapsed Time: 2.541 seconds (warm-up)
## Chain 2: 5.224 seconds (Sampling)
## Chain 2: 7.765 seconds (Total)

```

```

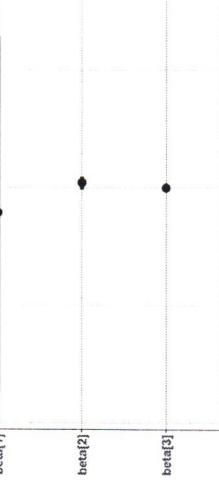
### Inference for Stan model: regression2
### 2 chains, each with iter=5000; warmup=1000; thin=5;
### post-warmup draws per chain=800, total post-warmup draws=1600.
###
```

	mean	se	sd	2.5%	25%	50%	75%	97.5%	n_eff	Rhat
beta[1]	-5.52	0.57	19.33	-43.46	-18.29	-5.38	6.90	32.05	1158	1
beta[2]	1.06	0.02	0.68	-0.29	0.69	1.05	1.50	2.37	1146	1
beta[3]	0.00	0.00	0.00	0.00	0.00	0.00	0.01	0.02	1137	1
sigma _{err}	200.87	6.65	244.93	158.35	194.97	200.85	217.72	255.73	21	1460

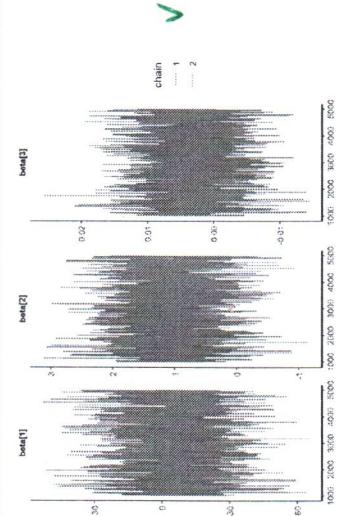
```
plot(fit3, ask = T, pars = c('beta1'), ci_level = 0.95, fill_color = "blue")
```

cl_level: 0.95
outer_level: 0.95 (95% intervals)

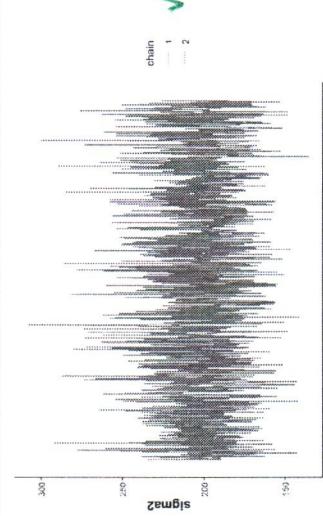
notaf1



```
# Draw the traceplot corresponding to one or more Markov chains,
# providing a visual way to inspect
# sampling behavior and assess mixing across chains and convergence
rstan::traceplot(fit3, pars = c("beta"), inc_warmup = FALSE)
```



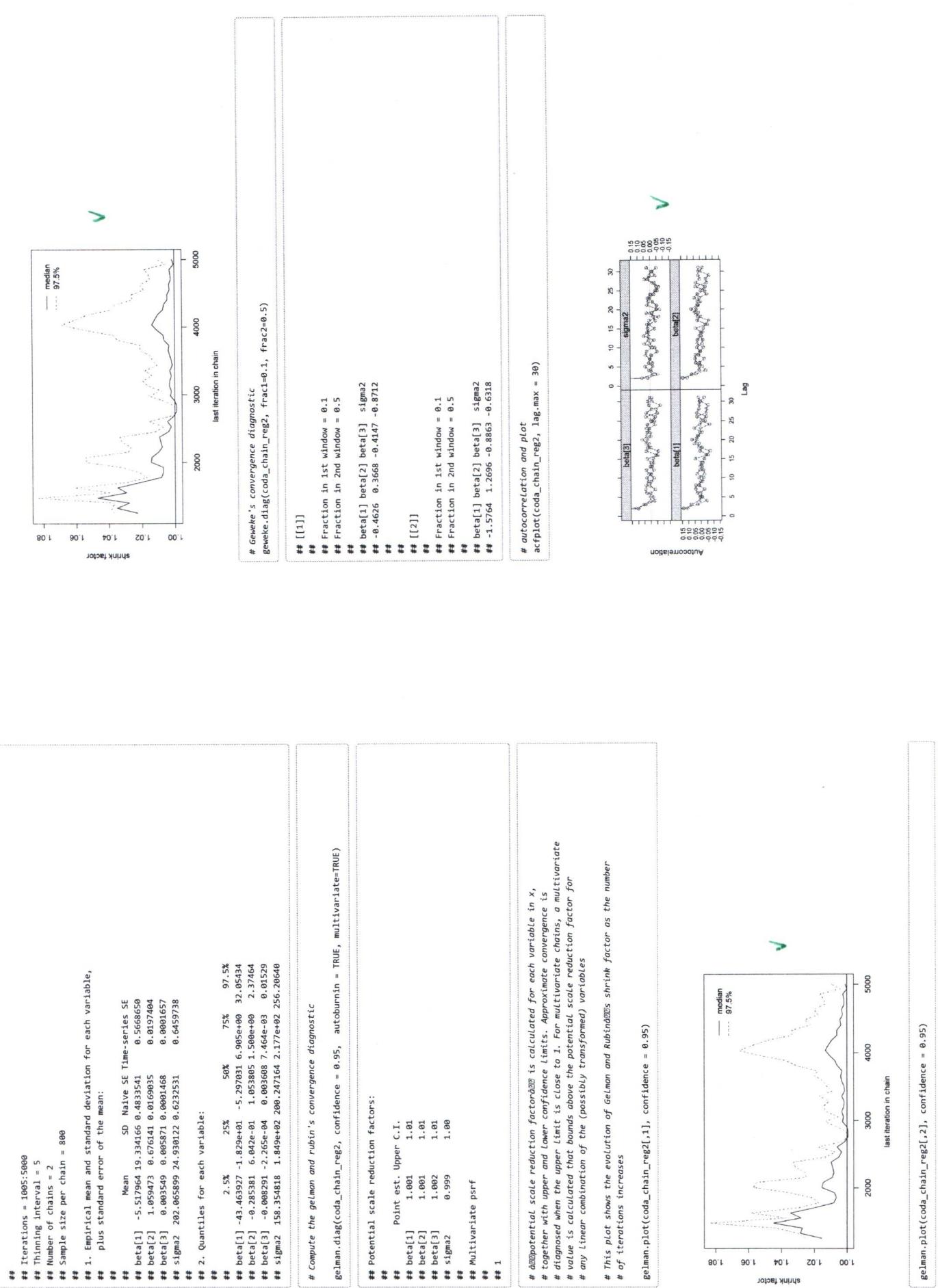
```
rstan::traceplot(fit3, pars='sigma2', inc_warmup = FALSE)
```



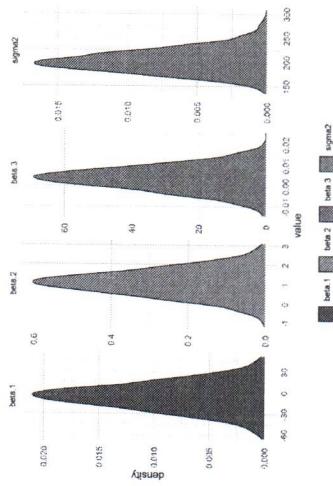
```

#-----#
# DIAGNOSTIC IN CODA
#-----#
coda_chain_reg <- as.mcmc.list(fit3, pars = c("beta", "sigma2"))
summary(coda_chain_reg)

```



```
# plots of the posteriors
plot_post <- fit3 %>%
  rstan::extract(c('beta', 'sigma2')) %>%
  map_df(as.data.frame, .id = 'param') +
  ggplot(aes(value, ffill = param)) +
  geom_density() +
  facet_wrap(~param, scales = 'free', nrow = 1) +
  scale_fill_manual() +
  theme_minimal() +
  theme(legend.position="bottom") +
  labs(fill = "")
```



```
## leave one out cross validation
loo(fit3)
```

leave one out cross validation

```
## Computed from 1600 by 133 log-likelihood matrix
## elpd_loo Estimate SE
## elpd_loo -546.5 8.6
## p_loo 2.1 0.4
## looic 1093.0 17.2
## Monte Carlo SE of elpd_loo is 0.0.
## All Pareto k estimates are good (k < 0.5).
## See help('pareto-k-diagnostic') for details.
```

loo(fit3)

lower vs higher

```
## Computed from 1600 by 133 log-likelihood matrix
## elpd_loo Estimate SE
## elpd_loo -45.0 9.3
## p_loo 3.9 0.7
## looic 1090.1 18.7
## Monte Carlo SE of elpd_loo is 0.1.
## All Pareto k estimates are good (k < 0.5).
## See help('pareto-k-diagnostic') for details.
```

```
# Look at Log-ic
# Lower is better, second model is preferred
```

```
my_WAIC <- function(fit, param){
  llk <- rstan:::extract(fit, param)[[1]]
  p_WAIC <- sum(sapply(llik, 2, var))
  lpd <- sum(sapply(llik, 2, function(x) log(mean(exp(x)))))
  WAIC <- -2 * lpd + 2 * p_WAIC
  return(WAIC)
}

my_WAIC(fit2, "log_lik")

## [1] 1093.025

my_WAIC(fit3, "log_lik")

## [1] 1090.068

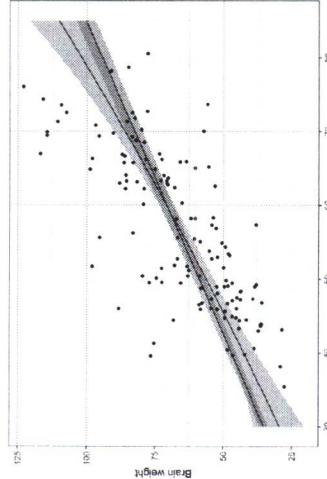
# asymptotically equivalent
# -----
# NICE PLOT
# -----
xseq <- seq(30, 85, length.out = 100)

df_temp <- t(apply(as.data.frame(rstan:::extract(fit2, pars = c('beta0', 'beta1'), permuted = TRUE)),
  1, function(x) x[1] + x[2] * xseq))

df_temp2 <- t(apply(as.data.frame(rstan:::extract(fit3, pars = c('beta0', 'beta1'), permuted = TRUE)),
  1, function(x) x[1] + x[2] * xseq + x[3] * xseq))

df_plot <- data.frame(x = xseq, y = colMeans(df_temp),
  ylow = apply(df_temp, 2, function(x) quantile(x, 0.025)),
  yup = apply(df_temp, 2, function(x) quantile(x, 0.975)),
  z = colMeans(df_temp2),
  zlow = apply(df_temp2, 2, function(x) quantile(x, 0.025)),
  yup = apply(df_temp2, 2, function(x) quantile(x, 0.975)))

ggplot(data.frame(x = regression_data$y, y = regression_data$y) +
  theme_bw() +
  geom_line(data = df_plot, aes(x = x, y = y), col = "blue") +
  geom_ribbon(data = df_plot, aes(x = x, ymin = ylow, ymax = yup), fill = "blue", alpha = 0.2) +
  geom_line(data = df_plot, aes(x = x, y = z), col = "red") +
  geom_ribbon(data = df_plot, aes(x = x, ymin = zlow, ymax = yup), fill = "red", alpha = 0.2) +
  geom_point(aes(x = x, y = y)) +
  xlab("SvL") +
  ylab("Brain weight")
```



*model + uncertainty
(in credible bounds)*

Variable selection and shrinkage

Riccardo Corradin
November 13, 2020

1/26

Introduction

The model choice problem

Strategy: fix **prior probability** masses to all the possible models, then apply the Bayes theorem to compute the **posterior distribution**.

The best models (more than one) are those with highest posterior probability.

Let m be the index describing L models M_1, \dots, M_L :

- $P(m = j)$, prior probability to choose model M_j (typically $1/L$)
- Model M_j :
 - likelihood $p(\mathbf{Y} | \theta_j, M_j)$
 - prior $\pi(\theta_j | M_j) = \pi(\theta_j | m = j)$
 - θ_j vector of parameters

We set a prior on each model and then we check the posterior
How can we do it?

We set a prior conditioned on the model (and the likelihood too)

Idea: scan the entire model space.

2/26

The model choice problem²

How it works...:

1. For any model M_j compute the posterior distribution of the parameter θ_j

$$\pi(\theta_j | \mathbf{Y}, M_j) = \frac{p(\mathbf{Y} | \theta_j, M_j) \pi(\theta_j | M_j)}{m(\mathbf{Y} | M_j)}$$

2. Compute the posterior probability masses for M_1, \dots, M_L

$$P(m = j | \mathbf{Y}) = \frac{m(\mathbf{Y} | M_j) P(m = j)}{m(\mathbf{Y})}, \quad j = 1, \dots, L$$

with

$$m(\mathbf{Y}) = \sum_{j=1}^L m(\mathbf{Y} | M_j) P(m = j)$$

3. Choose the model(s) with the highest $P(m = j | \mathbf{Y})$ and eventually compare these few models via predictive goodness-of-fit criteria

3/26

Problem of variable selection

In a regression world: given a dependent variable Y and a set of K potential regressors X_1, \dots, X_K , the problem is to find the **best** model of the form

$$g(\mathbb{E}[Y | \mathbf{X}, \boldsymbol{\beta}]) = \beta_1 X_1 + \dots + \beta_K X_K$$

i.e. linear models or generalized linear models, where we have at most K different covariates.

Using the previous exhaustive strategy, we have to estimate 2^K models to scan the entire space. With $K = 12$ we need to estimate 4096. The problem is intractable.

We need alternative strategies to deal with the model choice problem.

4/26

Spike and slab in a nutshell

A Hierarchical Mixture Model for Variable Selection

We index each of these possible 2^K subset choices by the vector

$$\boldsymbol{\gamma} = (\gamma_1, \dots, \gamma_K)^\top$$

where $\gamma_i = 0$ if $\beta_i = 0$, $\gamma_i = 1$ if $\beta_i \neq 0$, and $\{\gamma_s, s = 1, \dots, 2^K\}$ is the set of all possible models. We model the uncertainty underlying variable selection augmenting the problem with $\boldsymbol{\gamma}$, using a mixture prior

$$\pi(\boldsymbol{\beta}, \boldsymbol{\gamma}) = \pi(\boldsymbol{\beta} | \boldsymbol{\gamma})\pi(\boldsymbol{\gamma})$$

where

$$\gamma_j \sim \text{Bernoulli}(\theta_j), \quad j = 1, \dots, K, \quad \text{independent}$$

and θ_j is the probability that β_j is large enough to justify including X_j in the model.

5/26

Specifying the prior - Spike&Slab

This prior can also be used to put increased weight on parsimonious models by setting the θ_j small. The components of $\boldsymbol{\gamma}$ are apriori independent. Consider

$$\begin{aligned} \beta_j | \gamma_j &\sim (1 - \gamma_j)\delta_{\{0\}} + \gamma_j N(0, \sigma_{\beta_j}^2) \\ \gamma_j | \theta_j &\sim \text{Bern}(\theta_j) \\ \theta_j &\sim \pi(\theta_j) \end{aligned}$$

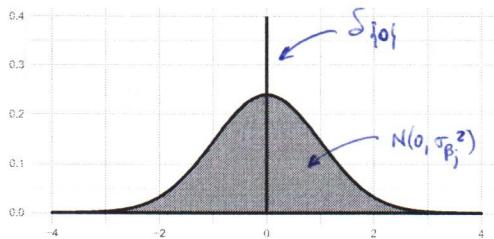
$\rightarrow \delta_j = 0 \Rightarrow$ all the mass in 0
 $\delta_j = 1 \Rightarrow$ all the mass in $N(0, \sigma_{\beta_j}^2)$

where $\delta_{\{0\}}$ is a Dirac measure with mass in 0 (you may also consider approximations of it!) and $N(0, \sigma_{\beta_j}^2)$ is a diffuse distribution.

A common choice for $\pi(\theta_j)$ is an $\text{Unif}(0, 1)$ distribution.

6/26

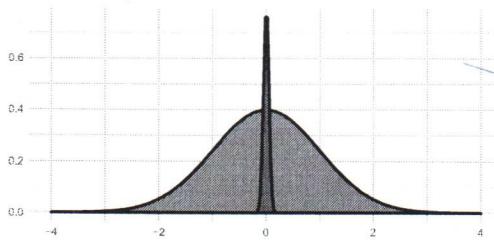
What are we doing? - Spike&Slab



Remark: in STAN is not possible to define a Spike&Slab prior we can use only continuous distributions for the parameters.

7/26

Relax the mixture - SSVS



for example we can set them both as gaussians

We can relax the mixture prior distribution on the coefficient by considering a combination of two continuous distribution: one super concentrated and the other slab.

8/26

Specifying the prior - SSVS

Consider a prior distribution for β_j as

$$\begin{aligned} \beta_j | \sigma_j^2 &\sim N(0, \sigma_j^2) \\ \sigma_j^2 | c_j, \tau_j, \gamma_j &\sim (1 - \gamma_j)\delta_{\tau_j^2} + \gamma_j \delta_{c_j^2 \tau_j^2} \\ \gamma_j | \theta_j &\sim \text{Bern}(\theta_j) \end{aligned}$$

→ if $\gamma_j = 1$ then the variance is $c_j^2 \tau_j^2$
otherwise it just τ_j^2
(c_j^2 = how much we amplify the variance in the slab component)

where $\delta_{\{x\}}$ is a Dirac measure with mass in x , and $N(0, \sigma_j^2)$ is a Gaussian (diffuse) distribution. We are combining together two diffuse distribution: a Gaussian with variance τ_j^2 and a Gaussian with variance $c_j^2 \tau_j^2$.

This prior looks like a relaxed version of the previous Spike&Slab prior.

9/26

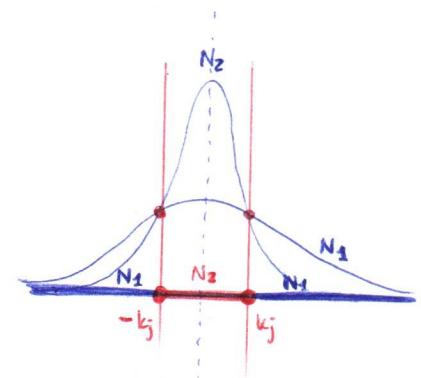
Specifying the prior - SSVS²

Within the previous priors setup we have a mixture of two Gaussian distributions, that intersect at the points $\pm \kappa_j = \tau_j \epsilon_j$ where

$$\epsilon_j = \sqrt{\frac{\log(c_j)c_j^2}{c_j^2 - 1}}$$

The point κ_j can be regarded as a threshold for declaring practical significance in that all coefficients falling into the interval $(-\kappa_j, \kappa_j)$ can be interpreted as 'practically zero'.

We use the previous priors in a practical example.



10/26

For both the cases (Spike & Slab / SSVS) we proceed similarly:

Model estimation

Let

- $\pi(\gamma)$: prior for the model index parameter γ
- $\pi(\beta | \gamma)$: prior for coefficients within each model
- $p(Y | \beta_\gamma, \gamma)$: likelihood (consider the product for the n subjects)

Model/covariate choice is made via the posterior of the model index parameter:

$$\pi(\gamma_j | Y) = \frac{p(Y | \gamma_j) \pi(\gamma_j)}{\sum_k p(Y | \gamma_k) \pi(\gamma_k)} \quad = \text{posterior distribution of including or not a variable in the model}$$

Where

$$\pi(Y | \gamma_j) = \int p(Y | \gamma_j, \beta_\gamma) \pi(d\beta_\gamma | \gamma_j)$$

is the marginal density of Y in a model indexed by γ_j .

11/26

Criteria for selecting covariates

Given (a sample from) the posterior distribution $\pi(\gamma | Y)$ we have different strategies to select the covariates:

- **HPD**: choose the model with the highest posterior probability

$$\hat{\pi}(\gamma_j | Y) = \frac{1}{n_{rep}} \sum_{t=1}^{n_{rep}} \mathbb{I}_{[\gamma_j^{(t)} = \gamma_j]} \quad = \text{we include those which prob. (posterior) is higher than a threshold}$$

- **MPM**: pick all variables with estimated marginal posterior inclusion probabilities

$$\hat{\pi}(\gamma_j | Y) = \frac{1}{n_{rep}} \sum_{t=1}^{n_{rep}} \mathbb{I}_{[\gamma_j^{(t)} = 1]}$$

- **HS**: hard shrinkage, pick all variables such that 0 does NOT belong to the marginal posterior **credible interval** of the corresponding regression parameter, or to a CI centered around the posterior mean $\pm \text{sd}$

12/26

Example:

Reach dataset

The Rotterdam Early Arthritis Cohort (dataset) study was initiated in 2004 to **investigate the development of rheumatoid arthritis** in patients with early manifestations of joint impairment.

Information regarding basic patient characteristics, serological measurements, and patterns of disease involvement at baseline has been gathered in **681 recruited patients**.

It is of interest to know which of the following **12 factors** are potentially **associated with the development** of rheumatoid arthritis considered as a binary (yes/no) outcome.

13/26

Covariates:

Reach dataset²

- ACCP (cyclic citrullinated peptide antibody)
- age
- ESR (erythrocyte sedimentation rate, is the rate at which red blood cells sediment in a period of one hour)
- DC (duration of complaints in days)
- stiffness (duration of morning stiffness in minutes)
- RF (rheumatoid factor)
- gender
- Sym (symmetrical pattern of joint inflammation; yes/no)
- SJC (swollen joint count)
- TJC (tender joint count)
- BCPH (bilateral compression pain in hands; yes/no)
- BCPF (bilateral compression pain in feet; yes/no)

14/26

Reach dataset³

The standard approach to analyze these data would be to use **probit regression** combined with some off-the-shelf variable selection method. The F-to-out backward selection with $p_D \leq 0.05$ yields a model with the following variables:

ACCP, ESR, DC, Sym, SJC, and BCPH

The model with the most favorable value of the AIC selected after an exhaustive model evaluation contains two extra variables: RF and stiffness.

We consider a Bayesian probit regression, where the coefficients are selected via Spike&Slab priors.

$$\begin{aligned} \text{probit}(\mathbb{E}[Y_i | \mathbf{X}_i]) &= \mathbf{X}_i^T \boldsymbol{\beta} \\ \text{probit}(a) = \Phi^{-1}(a) \implies \mathbb{E}[Y_i | \mathbf{X}_i] &= \Phi(\mathbf{X}_i^T \boldsymbol{\beta}) \\ \Phi(\cdot) \text{ denote the CDF of a Gaussian distribution} \end{aligned}$$

15/26

Recap on the model

Model specification:

$$\begin{aligned} Y_i | \mathbf{X}_i, \boldsymbol{\beta} &\sim \text{Bern}(\Phi(\mathbf{X}_i^T \boldsymbol{\beta})), \quad i = 1, \dots, n \\ \boldsymbol{\beta}_{\gamma} | \gamma &\sim \dots \quad \xleftarrow{\text{depending on if we decide to}} \\ \gamma_j | \theta_j &\sim \text{Bern}(\theta_j), \quad j = 1, \dots, K \end{aligned}$$

depending on if we decide to consider Spike & Slab or SSVS

Prior elicitation:

- we set $\boldsymbol{b} = 0$ (vector), expected value for the distribution on $\boldsymbol{\beta}$, i.e. considering the symmetry of the distribution we are not saying that a priori the effect is positive or negative.
- working with marginally standardized covariates
- eventual hyperpriors to elicitate the model

16/26

See the code!

Regularization

Brief introduction to regularization methods

Data: $\{(X_1, Y_1), \dots, (X_n, Y_n)\}$. Model specification:

$$g(\mathbb{E}[Y_i | X_i]) = \mathbf{X}_i^T \boldsymbol{\beta}$$

Problem: the regression coefficients may explode, the best fit can be poor in terms of interpretability.

Maximum likelihood estimation:

$$\hat{\boldsymbol{\beta}}_{ML} = \arg \max_{\boldsymbol{\beta}} \left\{ \sum_{i=1}^n \log f(Y_i | \mathbf{X}_i, \boldsymbol{\beta}) \right\}$$

Penalized estimation:

$$\hat{\boldsymbol{\beta}}_{penalized} = \arg \max_{\boldsymbol{\beta}} \left\{ \sum_{i=1}^n \log f(Y_i | \mathbf{X}_i, \boldsymbol{\beta}) + h(\boldsymbol{\beta}) \right\}$$

by including this term we can push the value of the parameters to zero

17/26

Brief introduction to regularization methods²

Different choices for the shrinkage penalty lead to different regularization methods:

- A L_1 norm leads to the **Lasso** method:

$$h(\boldsymbol{\beta}) = \lambda \|\boldsymbol{\beta}\|_1 = \lambda \sum_{j=1}^K |\beta_j|$$

- A L_2 norm leads to the **Ridge** method:

$$h(\boldsymbol{\beta}) = \lambda \|\boldsymbol{\beta}\|_2 = \lambda \sum_{j=1}^K \beta_j^2$$

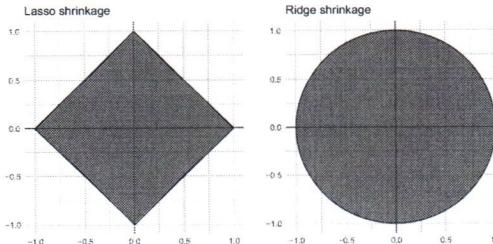
- Combining together L_1 and L_2 norms we get the **Elastic Net**:

$$h(\boldsymbol{\beta}) = \lambda_1 \|\boldsymbol{\beta}\|_1 + \lambda_2 \|\boldsymbol{\beta}\|_2 = \lambda_1 \sum_{j=1}^K |\beta_j| + \lambda_2 \sum_{j=1}^K \beta_j^2$$

18/26

Brief introduction to regularization methods³

But what are we doing practically?



With the **Lasso** shrinkage we are forcing the parameter to be inside a square, with the **Ridge** inside a circle. With the **Elastic Net** we are combining the two penalties.

19/26

Brief introduction to regularization methods⁴

The **Lasso** penalty is known for its sparsity property, i.e. is forcing some coefficients quickly to zero.

The **Ridge** penalty is known for its grouping property, in the sense that it is able to retain groups of strongly correlated variables.

The **Elastic Net** is a good compromise between this two properties.

But what happens in **Bayesian** statistics?

20/26

Regularization methods with Bayes

The choice of β which minimize the **Ridge** penalty is also the value which maximize the probability of independent random variables distributed as a **Gaussian distribution**. Let assume that $\beta | \lambda \sim N(0, I_K / \lambda)$.

$$\begin{aligned}\tilde{\beta} &= \arg \max_{\beta} \left\{ (2\pi)^{-K/2} \lambda^{K/2} e^{-\frac{1}{2} \beta^T \beta} \right\} \\ &= \arg \max_{\beta} \left\{ -\frac{\lambda}{2} \beta^T \beta \right\} \\ &= \arg \min_{\beta} \left\{ \lambda \sum_{j=1}^K \beta_j^2 \right\}\end{aligned}$$

where λ is a common penalty.

21/26

Regularization methods with Bayes²

Thus the **Ridge** regression correspond to the maximum a posteriori (MAP) estimation of a model with a Gaussian distribution on the regression coefficients and a particular structure for the variance.

The specification for the model is then

$$\begin{aligned}Y_i | \mathbf{X}_i, \beta &\sim p(Y_i | \mathbf{X}_i, \beta), \quad i = 1, \dots, n \\ \beta_j | \lambda &\sim N(0, 1/\lambda), \quad j = 1, \dots, K\end{aligned}$$

About the parameter λ :

- we can fix the parameter to a specific value
- we can specify an hyperprior distribution, for example

$$\lambda \sim \text{Gamma}(a_\lambda, b_\lambda)$$

22/26

Regularization methods with Bayes³

The choice of β which minimize the **Lasso** penalty is also the value which maximize the probability of independent random variables distributed as a **Laplace distribution**. Let assume that $\beta_j | \lambda \sim \text{Lap}(0, 1/\lambda_j)$.

$$\begin{aligned}\tilde{\beta} &= \arg \max_{\beta} \left\{ \prod_{j=1}^K \frac{\lambda_j}{2} e^{-\lambda_j |\beta_j|} \right\} \\ &= \arg \max_{\beta} \left\{ \sum_{j=1}^K -\lambda_j |\beta_j| \right\} \\ &= \arg \min_{\beta} \left\{ \sum_{j=1}^K \lambda_j |\beta_j| \right\} \\ &= \arg \min_{\beta} \left\{ \lambda \sum_{j=1}^K |\beta_j| \right\}\end{aligned}$$

where the last holds assuming a common penalty.

23/26

Regularization methods with Bayes⁴

Thus the **Lasso** regression correspond to the maximum a posteriori (MAP) estimation of a model with a Laplace distribution on the regression coefficients.

The specification for the model is then

$$\begin{aligned}Y_i | \mathbf{X}_i, \beta &\sim p(Y_i | \mathbf{X}_i, \beta), \quad i = 1, \dots, n \\ \beta_j | \lambda &\sim \text{Lap}(0, 1/\lambda), \quad j = 1, \dots, K\end{aligned}$$

About the parameter λ :

- we can fix the parameter to a specific value
- we can specify an hyperprior distribution, Park and Casella suggested

$$\lambda^2 \sim \text{Gamma}(a_\lambda, b_\lambda)$$

24/26

Regularization methods with Bayes³

In a similar spirit, the **Elastic Net** in the Bayesian field correspond to a particular choice of the prior distribution on the regression coefficients, such that the MAP estimation of the parameters is equivalent to the ones that minimize the **Elastic Net** penalty.

The prior distribution for β is an hierarchical one, defined as:

$$\begin{aligned}\beta_j \mid \tau_j, \sigma^2 &\sim N\left(0, \left[\frac{\lambda_2}{\sigma^2} \frac{\tau_j}{\tau_j - 1}\right]^{-1}\right) \\ \tau_j \mid \sigma^2 &\sim \text{trGamma}\left(0.5, \frac{\lambda_1^2}{8\lambda_2\sigma^2}, 1, \infty\right)\end{aligned}$$

where $\text{trGamma}(a, b, c, d)$ correspond to the Gamma distribution with shape parameter a and rate parameter b , truncated on (c, d) .

25/26

Regularization methods with Bayes²

The **Elastic Net** correspond to the maximum a posteriori (MAP) estimation of a model with the previous prior specification on the regression coefficients.

The specification for the model is then

$$\begin{aligned}Y_i \mid \mathbf{X}_i, \beta &\sim p(Y_i \mid \mathbf{X}_i, \beta), \quad i = 1, \dots, n \\ \dots &\sim \dots\end{aligned}$$

plus the previous prior for the coefficients β . About the parameter λ_1, λ_2 :

- we can fix the parameter to a specific value
- we can specify an hyperprior distribution, Park and Casella suggested

$$\lambda_j^2 \sim \text{Gamma}(a_{\lambda_j}, b_{\lambda_j}), \quad j = 1, 2$$

26/26

See the code!

```

## BL_probit2.stan
# //#####
# //### BAYESIAN ELASTIC NET: PROBIT REGRESSION #####
# //#####
# //#####
# //##### DATA //#####
# //#####
# data {
#   # int<lower = 0> N; // number of data
#   # int<lower = 0> p; // number of covariates (without the intercept)
#   # int<lower = 0> upper = 1> Y[N]; // response vector
#   # matrix[N, p] X; // design matrix
#   # real<lower = 0> Lambda;
# }

# //#####
# parameters {
#   vector[p] beta; // regression coefficients
#   real betao; // intercept
# }

# //#####
# model {
#   # //Likelihood
#   # for (s in 1:N)
#   # {
#   #   // print(" s = ", s);
#   #   Y[s] ~ bernoulli(phi(beta0 + row(X, s)*beta));
#   # }

#   # // Prior
#   # // beta
#   # betao ~ normal(0, 50.0); // vague prior: (mean, sd)
#   # for (j in 1:p)
#   # {
#   #   beta[j] ~ double_exponential(theta, 1.0 / Lambda);
#   # }
# }

## BR_probit2.stan
# //#####
# //### BAYESIAN ELASTIC NET: PROBIT REGRESSION #####
# //#####
# //#####
# //##### DATA //#####
# //#####
# data {
#   # int<lower = 0> N; // number of data
#   # int<lower = 0> p; // number of covariates (without the intercept)
#   # int<lower = 0> upper = 1> Y[N]; // response vector
#   # matrix[N, p] X; // design matrix
#   # real<lower = 0> Lambda;
# }

# //#####
# parameters {
#   vector[p] beta; // regression coefficients
#   real betao; // intercept
# }

# //#####
# model {
#   # //Likelihood
#   # for (s in 1:N)
#   # {
#   #   // print(" s = ", s);
#   #   Y[s] ~ bernoulli(phi(beta0 + row(X, s)*beta));
#   # }

#   # // Prior
#   # // beta
#   # betao ~ normal(0, 0.5); // vague prior: (mean, sd)
#   # for (j in 1:p)
#   # {
#   #   beta[j] ~ normal(theta, pow(1.0 / Lambda, 0.5));
#   # }
# }

```

```
data_JAGS_SpS1 <- list(N = N, p = p, Y = Y, X = as.matrix(X),
var_beta = rep(1, p))
```

```
#####
##### VARIABLE SELECTION #####
#####

# R & STAN are friends!
library(doParallel)
library(rstan)
library(rjags)
library(coda)
library(ggplot2)
library(dplyr)
library(tidyverse)
library(purrr)
library(BGmisc)
require(ggplots)
require(ggpubr)

#####
##### REACH data #####
## Rheumatoid arthritis is an autoimmune disease characterized by
## chronic synovial inflammation and destruction of cartilage and bone
## In the Joints.
# The Rotterdam Early Arthritis Cohort (REACH) study was
## initiated in 2004 to investigate the development of rheumatoid arthritis
## in patients with early manifestations of joint impairment (compromissione articolare).
## Information regarding basic patient characteristics, serological
## measurements, and patterns of disease involvement at baseline has
## been gathered in 681 recruited patients. It is of interest
## to know which of the following 12 factors are potentially
## associated with the development of rheumatoid arthritis considered as
## a binary (yes/no) outcome:
## ACP (cyclic citrullinated peptide antibody),
## age,
## ESR (erythrocyte sedimentation rate, is the rate at which red blood cells sediment in a period of one
## hour),
## DC (duration of complaints in days),
## stiffness (duration of morning stiffness in minutes),
## RF (rheumatoid factor),
## gender,
## TJC (swollen joint count),
## SJC (tender joint count),
## BCPH (bilateral compression pain in hands; yes/no),
## ECPF (bilateral compression pain in feet; yes/no),
## The standard approach to analyze these data would be to use
## logistic/probit regression combined with some off-the-shelf variable
## selection method. The F-to-out backward selection with p_D 0.05 yields a
## model with the following variables:
## -----> ACP, ESR, DC, Sym, SJC, and BCPH.
## The model with the most favorable value of the AIC selected after an
## exhaustive model evaluation contains two extra variables: RF and stiffness.
## which of these models provide the best approximation to the true
## underlying relationships is, if at all possible, difficult to assess.
```

reach <- read.table("REACH_data.txt", header=TRUE)

names(reach)

head(reach)

reach\$gender <- reach\$gender - 1

sub.idv = 1:12

num.data = 681

X <- as.matrix(reach[1:num.data, sub.idv])

Y <- as.vector(reach[1:num.data, 13])

N <- dim(X)[1]

p <- dim(X)[2]

#####
Spike and Slab
#####

data to pass to JAGS (see the code in SSVS_probit.bug)

```
#####
##### A list of initial value for the MCMC algorithm
## that WinBUGS will implement.
inits = function() {
  list(bao = 0.0, btemp = rep(0,p), g = rep(0,p), theta = rep(0.5, p),
.RNG.seed = 321, .RNG.name = 'base: Wichmann-Hill')
}

model=jags.model("SpS1_probit.bug",
data = data_JAGS_SpS1,
n.adapt = 1000,
inits = inits,
n.chains = 1)

# if we want to perform a larger burn in with not adaptation.
#update(model,n.iter=500)

# Posterior parameters WinBUGS has to save NO "sigma2" here
param <- c("betaad", "beta", "g", "mell")

# number of iterations & thinning
nit <- 50000
thin <- 10

## The command coda.sample() calls jags from R passing the data and initial value just defined
output <- coda.samples(model = model,
variable.names = param,
n.iter = nit,
thin = thin)

# save(output,file='SpS1.dat') # we save the chain
load('SpS1.dat')

## We give a look at the trace plot and at the density summary
## of the posterior chains for each of the parameters
str(output)

## The output is an mcmc object of the library coda
plot(output,ask=T)

## We give a look at the trace plot and at the density summary
## of the posterior chains for each of the parameters
str(output)

## Some variable selection techniques:
## The median probability model (MPM)
## pick variables with estimated posterior inclusion probabilities
## higher than 0.5
## Notice that the estimated posterior inclusion probabilities are the
## posterior means of the Gamma variables (in the code we called g)
head(output)

#####
##### save the posterior chain of the inclusion variable in post_g
post_g <- as.matrix(output[,14:25])
apply(post_g,2,"mean")
post_mean_B <- apply(post_g,2,"mean")

#####
##### plot the posterior inclusion probabilities
p1 <- data.frame(value = post_mean_E, var = colnames(X)) %>%
ggplot(aes(y = value, x = var, fill = var)) +
geom_bar(stat="identity") +
geom_line(mapping = aes(intercept = .5), col = 2, lwd = 1.1) +
coord_flip() +
theme_minimal() +
theme(legend.position="none") +
ylab("Posterior inclusion probabilities") +
xlab("")
```

p1

→ posterior probability of including a specific covariate in the model
(we tune/knowsider 30.5)

we will compare later the model with other methods

mp_SpS1 <- as.vector(which(post_mean_E > 0.5))

```

post_mean_g[mp_SSV1]

# ANOTHER WAY TO CHOOSE THE MODEL:
# Highest posterior density model (HPD)
# Pick a model with the highest estimated posterior probability

# Recall that we have represented the model index
# using a binary coding as
# md[i]=2^i+2^(i-1)+...+2^0BP
# the visited model are saved in the chain
plot(output[, "md1"], pch = 20)

# for example at iteration 10 the chain explored the
# model
output[10, "md1"]

##### We start to analyze how many models have been visited
# by the posterior chain:
length(unique(output[, "md1"]))

##### Now we compute the posterior frequency of the visited models
# model visited out of 4096
# and sort the results
unique_model <- unique(post_B, MARGIN = 1)
freq <- apply(unique_model, 1, function(b) sum(apply(post_E, MARGIN = 1, function(a) all(a == b))))
cbind(unique_model[order(freq,decreasing = T)], sort(freq,decreasing = T))

##### the HPD model is
colnames(X)[as.logical(unique_model[[which.max(freq)]])]

HPD_SPSI <- c(1:12)/as.logical(unique_model[[which.max(freq)]])]

#####
##### SSSVS #####
#####

##### Parameters of the spike slab prior Set 1
c_ss <- 100
intersect <- intersect / sqrt(2 * log(c_ss) * c_ss^2/(c_ss^2 - 1))

##### With this choice of hyperparameter, c_ss and intersection
##### the variance of the quasi-spike prior is
tau_ss^2

##### While the variance of the slab is
(tau_ss*c_ss)^2

##### A list of initial value for the MCMC algorithm
# the spike is super concentrated
curve(dnorm(x, 0, tau_ss), xlim = c(-5, 5))
curve(dnorm(x, 0, tau_ss * c_ss), xlim = c(-5, 5), add = T, col = 2)
# that WinBUGS will implement
initS = function() {
  list(beta0 = 0.0, beta = rep(0,p), g = rep(0,p),
       tau_ss = tau_ss, c_ss = c_ss)
}

##### data to pass to JAGS (see the code in SSVS_probbit.bug)
data_JAGS_1 <- list(N = N, p = p, Y = Y, X = as.matrix(X),
                     tau_ss = tau_ss, c_ss = c_ss)

##### A list of initial value for the MCMC algorithm
# that WinBUGS will implement
initS = function() {
  list(beta0 = 0.0, beta = rep(0,p), g = rep(0,p),
       RNG.seed = 321, .RNG.name = 'base: Wichmann-Hill')
}

model=jags.model("SSVS_probbit.bug",
                 data = data_JAGS_1,
                 n.adapt = 10000,
                 inits = initS,
                 n.chains = 1)

# If we want to perform a larger burn in with not adaptation.

# Another way to choose the model:
# Highest posterior density model (HPD)
# pick a model with the highest estimated posterior probability
# Recall that we have represented the model index
# using a binary coding as
# md[1:2^p]+...+2^0BP
# the visited model are saved in the chain
plot(output[, "md1"], pch = 20)

# for example at iteration 10 the chain explored the
# model

```

```

## we start to analyze how many models have been visited
## by the posterior chain:
length(unique(output[, "md1"]))

# model visited out of 4996
## Now we compute the posterior frequency of the visited models
visited_models<-table(output[, "md1"])

visited_models

# HPD model
# getting the unique profiles
# and sort the results
unique_model <- unique(post_B, margin = 1)
freq <- apply(unique_model, 1, function(b) sum(apply(post_B, margin = 1, function(a) all(a == b))))
cbind(unique_model[order(freq, decreasing = T),], sort(freq, decreasing = T))

# the HPD model is
colnames(X)[as.logical(unique_model[which.max(freq),])]

HP_SSV1 <- c(1:12)[as.logical(unique_model[which.max(freq),])]

#----- SCENARIO 2 -----
# Parameters of the spike slab prior set 2
c_ss < 100
intersect <- 0.1
tau_ss <- intersect / sqrt(2 * log(c_ss) * c_ss^2 / (c_ss^2 - 1))

## With this choice of hyperparameter, c_ss and intersection
## the variance of the quasi-spike prior is
tau_ss2

#while the variance of the slab is
(tau_ss*c_ss)^2

#---
curve(dnorm(x, 0, tau_ss), xlim = c(-5, 5),
      curve(dnorm(x, 0, tau_ss * c_ss), xlim = c(-5, 5), add = T, col = 2)
      # the spike is super concentrated
      # that WinBUGS will implement
      inits2 = function() {
        list(beta0 = rep(0, p), beta = rep(0, p), g = rep(0, p),
             n.adapt = 1000,
             inits = inits2,
             n.chains = 1)
      }

      # A list of initial value for the MCMC algorithm
      # that WinBUGS will implement
      inits2 = function() {
        list(beta0 = rep(0, p), beta = rep(0, p), g = rep(0, p),
             n.adapt = 1000,
             inits = inits2,
             n.chains = 1)
      }

      # if we want to perform a larger burn in with not adaptation.
      update(model,n.iters=500)

      # Posterior parameters WinBUGS has to save NO "sigma2" here
      param <- c("beta0", "beta", "g", "md1")

      # Some other information for the
      # MCMC algorithm
      # number of iterations
      n.iters <- 50000

      # thinning
      thin <- 10

      ##The command code.sample() calls jags from R passing the data and initial value just defined
      output <- code.sample(model = model,
                           variable.names = param,
                           n.iters = n.iters,
                           thin = thin)

      # shrinkage
      strParam <- rstan::extract(BL_model, pars = c("lambda2"))
      permuted = TRUE
      rstan::traceplot(BIN_model, pars = c("lambda2"))

      load('BL_model.dat')

      is(BL_model)
      names(BL_model)
      summary(BL_model)

      # shrinkage
      strParam <- rstan::extract(BL_model, pars = c("lambda2"))
      permuted = TRUE
      rstan::traceplot(BIN_model, pars = c("lambda2"))
    )
  )
}

```

```

error = error + abs(Y[1] - ifelse(mean(out) > 0.5, 1, 0))
segments[, Y[1], col = 'royalblue', lwd = 2)
points[, Y[1], col = 'darkred', pch = 19)
points[, mean(out), col = 'royalblue', pch = 19)
message("### Subject ", i)

}
abline(h = 0.5)
errBl <- error/subN

#####
##### BAYESIAN ELASTIC NET #####
#####

data.win <- list(N = N, p = p, Y = Y, X = as.matrix(X))
inits <- function()
{
  list(beta0 = 0, beta = rep(0,1,p), tau = rep(2, p), a1 = 10, a2 = 20)
}

#####
##### FIT THE MODEL #####
#####

# If the credibility interval does not contain 0
# then I keep the variable
if(CL_beta <= apply(beta, 2, quantile, c(0.025, 0.975)))
CL_beta = apply(beta, 2, quantile, c(0.025, 0.975))

# Extract values of the chain
par <- rstan::extract(BE_model, pars = c("beta0", "beta"), permuted = TRUE)

#####
##### Variable selection:
# compute the 95% posterior credible interval for beta
CL_beta = apply(beta, 2, quantile, c(0.025, 0.975))

length(beta0)
beta <- as.matrix(par$beta)
dim(beta)

# If the credibility interval does not contain 0
# then I keep the variable
if(CL_beta[1,1]<0 && CL_beta[2,1]>0)
{
  cat("### variable ", colnames(reach)[1], " excluded \n")
}
else
{
  cat("### variable ", colnames(reach)[1], " included \n")
  idx_cov_BI = c(idx_cov_BI, 1)
}
}

mean_beta_post <- apply(beta, 2, "mean")
mean_beta_post

#####
##### boxplot #####
# boxplot
data.frame(x = as.vector(beta), var = rep(colnames(reach)[1:12], each = nrow(beta))) %>%
ggplot(aes(x = var, y = x, fill = var)) +
geom_boxplot() +
theme_bw() +
theme(legend.position = "null") +
geom_hline(aes(yintercept = 0), col = 2, lty = 2) +
theme(panel.spacing = element_text(angle = 45, hjust = 1)) +
xlab("") +
ylab("") +
yaxis.title = element_text(size = 10, color = "black", family = "sans-serif")

#####
##### PREDICTION in the PROBABILITY SCALE:
# loop over the iterations
out <- numeric(S)
for(g in 1:G) out[g] <- rnorm(sum(Xnew$beta[g]), mean = 0, sd = 1)

data.frame(out = out) %>%
ggplot(aes(out)) +
geom_histogram(bins = 40, alpha = 0.6, col = 1) +
theme_minimal() +
theme(legend.position="none") +
geom_vline(mapping = aes(xintercept = 1), col = 2, lwd = 1.2) +
xlab("") +
yaxis.title = element_text(size = 10, color = "black", family = "sans-serif")

#####
##### EVALUATE THE PREDICTION ERROR ON THE OBSERVED SAMPLE:
error = 0
subN = 100

for(g in 1:subN){
  Xnew = X[1,]
  for(g in 1:G) # loop over the iterations
  out[g] = rnorm(sum(Xnew$beta[g]), mean = 0, sd = 1)

  plot(1:subN, rep(0.5, subN), col = "white", ylim = c(-0.01, 1.01), xlab = "", ylab = "")
  for(g in 1:subN){
    Xnew = X[1,]
    for(g in 1:G) # loop over the iterations
    out[g] = rnorm(sum(Xnew$beta[g]), mean = 0, sd = 1)
  }
}

#####
##### Posterior predictive distribution #####
#####

# WARNING: the model is a bit slow to run...
BE_model = stan(file = "BEN_prob0.stan", data = data.win,
chains = 1, iter = 25000, warmup = 5000, thin = 5, seed = 199,
init = inits, control = list(max_treedepth = 20, adapt_delta = 0.95))
# save(BE_model, file="BEN_model.dat")
load("BEN_model.dat")

is(BE_model)
names(BE_model)
summary(BE_model)

# shrinkage
shr_param <- rstan::extract(BE_model, pars = c("a1", "a2"),
rstan::traceplot(BE_model, pars = c("a1", "a2")))

plot_post <- shr_param %>%
as_tibble() %>%
map_df(as.data.frame, .id = 'param')

plot_post %>%
ggplot(aes(value, fill = param)) +
geom_density() +
facet_wrap(~param, scales = 'free') +
scale_fill_manual() +
theme_minimal() +
theme(legend.position="bottom")

# Extract values of the chain
par <- rstan::extract(BE_model, pars = c("beta0", "beta"), permuted = TRUE)
names(par)
beta0 <- par$beta0;
length(beta0)
beta <- as.matrix(par$beta);
dim(beta)

#####
##### Variable selection:
# compute the 95% posterior credible interval for beta
CL_beta = apply(beta, 2, quantile, c(0.025, 0.975))
CL_beta

# If the credibility interval does not contain 0
# then I keep the variable
idx_cov_BEN = NULL
for(l in 1:p){
  if(CL_beta[l,1]<0 && CL_beta[l,2]>0)
  {
    cat("### variable ", colnames(reach)[l], " excluded \n")
  }
  else
  {
    cat("### variable ", colnames(reach)[l], " included \n")
    idx_cov_BEN = c(idx_cov_BEN, 1)
  }
}

#####
##### Covariance matrix #####
#####

Xnew <- X[1,]
for(g in 1:G) Xnew = rnorm(sum(Xnew$beta[g]), mean = 0, sd = 1)

# loop over the iterations
out <- numeric(S)
for(g in 1:G) out[g] <- rnorm(sum(Xnew$beta[g]), mean = 0, sd = 1)

data.frame(out = out) %>%
ggplot(aes(out)) +
geom_histogram(bins = 40, alpha = 0.6, col = 1) +
theme_minimal() +
theme(legend.position="none") +
geom_vline(mapping = aes(xintercept = 1), col = 2, lwd = 1.2) +
xlab("") +
yaxis.title = element_text(size = 10, color = "black", family = "sans-serif")

```

```

}
}

mean_beta_post <- apply(beta, 2, "mean")
mean_beta_post

#-----#
* boxplot
data.frame(x = as.vector(beta), var = rep(colnames(reach)[1:12], each = nrow(beta))) %>%
  ggplot(aes(x = var, y = x, fill1 = var)) +
  geom_boxplot() +
  theme_bw() +
  theme(legend.position = "null") +
  theme_hline(aes(y intercept = 0), col = 2, lty = 2) +
  theme(axis.text.x = element_text(angle = 45, hjust = 1)) +
  xlab("") +
  ylab("")

# NOTE: I can also choose a different interval (80%, 90%)
# ----> the smaller the degree is, the less I am restrictive in choosing covariates

## PREDICTION in the PROBABILITY SCALE:
Xnew <- X[,]
G <- 4600 # num iteration
out <- numeric(G)
for (g in 1:G) out[g] <- rnorm(sum(Xnew*beta[,]), mean = 0, sd = 1)

data.frame(out = out) %>%
  ggplot(aes(out)) +
  geom_histogram(bins = 40, alpha = 0.6, col = 1) +
  geom_vline(mapping = aes(xintercept = 1), col = 2, lwd = 1.1) +
  theme_minimal() +
  theme(legend.position="none") +
  ylab("predicted probability") +
  xlab("")

#-----# EVALUATE THE PREDICTION ERROR ON THE OBSERVED SAMPLE:
error = 0
subN = 100
registerParallel(cores=no_cores)

model_list <- foreach(i = 1:length(lambda_seq)) %dopar% {
  Xnew = X[,]
  for(g in 1:subN){
    out[g] = prnorm(sum(Xnew*beta[,]), mean = 0, sd = 1)
    error + abs((out[i] - ifelse(out > 0.5, 1, 0)))
  }
  segments(i, quantile(out, 0.95), i, quantile(out, 0.05), col = "royalblue", lwd = 2)
  points(i, Y11, col = 'darkred', pch = 19)
  points(i, mean(out), col = 'royalblue', pch = 19)
  message("*** Subject: ", i)
}
errBEN <- error/subN

#-----# compare
pB1 <- ggplot(df1) +
  geom_line(aes(x = lambda, y = Y, X = as.matrix(X), lambda = lambda_seq[1])) +
  theme_bw() +
  ylab("posterior mean") +
  scale_color_discrete(name = "variable") +
  ggtitle("BAYESIAN RIDGE")

pB2 <- ggplot(df2) +
  geom_line(aes(x = lambda, y = Y, X = as.matrix(X), lambda = lambda_seq[1])) +
  theme_bw() +
  ylab("posterior mean") +
  scale_color_discrete(name = "variable") +
  ggtitle("BAYESIAN LASSO")

#-----# errors
errBL <- errBEN
errBN <- errBEN

#-----# Bayesian ridge
mp_SpS1 # median probability model (SpS1)
HDP_SpS1 # Highest posterior density model (HDP) (SpS1)
mp_SSV1 # median probability model (SSV1)
HDP_SSV1 # Highest posterior density model (HDP) (SSV1)
mp_SSV2 # median probability model (SSV2)
HDP_SSV2 # Highest posterior density model (HDP) (SSV2)
idx_cov_BL # Bayesian lasso
idx_cov_BN # Bayesian elastic net

#-----# Bayesian Lasso
mp_model <- stan(file = 'Bl_probit2.stan', model_name = 'Bl_stan')
compiled_model <- stan_model(stanc.ret = build_model, verbose = FALSE)

```

Mixing the models and Time series data

Riccardo Corradin
November 19, 2020

1/20

Introduction

GLM

Let Y be the response variable, $Y \in \mathcal{Y}$, and x a p -dimensional vector of covariates, with $x \in \mathcal{X}$. Assume that $f_Y \in \mathcal{P}$ exponential family of distribution.

Moreover, let $\theta = (\beta, \xi) \in \Theta$ a vector of parameters, with Θ parameter space, where β denotes the regression coefficients and ξ denotes the other parameters, such that the following holds

$$\mathbb{E}[Y | x, \beta, \xi] = g^{-1}(x^\top \beta)$$

A Generalized Linear Model is composed by three terms:

- A probability distribution f_Y (belonging to the exponential family)
- A linear predictor $\eta = x^\top \beta$
- A link function $g(\cdot)$ such that $\mathbb{E}[Y | x, \beta, \xi] = g^{-1}(\eta)$

2/20

GLM²

Different probability distributions and link functions lead to different models. Denote by $\mu = \mathbb{E}[Y | x, \beta, \xi]$. Examples of models:

- Normal distribution for Y
 - Identity link function, leads to the linear regression

$$x^\top \beta = \eta = g(\mu) = \mu$$

- Bernoulli distribution for Y (or categorical)
 - Logit link function, leads to the logistic regression

$$x^\top \beta = \eta = g(\mu) = \log\left(\frac{\mu}{1-\mu}\right)$$

- Probit link function, leads to the probit regression

$$x^\top \beta = \eta = g(\mu) = \Phi^{-1}(\mu)$$

- Poisson distribution for Y
 - Log link function, leads to the Poisson regression

$$x^\top \beta = \eta = g(\mu) = \log(\mu)$$

3/20

GLM with Bayes

Assume that we have an exchangeable sample $\{(Y_1, \mathbf{x}_1), \dots, (Y_m, \mathbf{x}_n)\}$. We denote by

$$\mathcal{L}(\mathbf{Y} | \mathbf{X}, \boldsymbol{\theta}) = \prod_{i=1}^n f_Y(Y_i | \mathbf{x}_i, \boldsymbol{\theta})$$

the **likelihood** of our data, where $f_Y(Y_i | \mathbf{x}_i, \boldsymbol{\theta})$ corresponds to the probability distribution of a **GLM** model, and by $\pi(\boldsymbol{\theta})$ the prior distribution on the vector of parameters $\boldsymbol{\theta}$.

We are interested into performing inference on the parameters $\boldsymbol{\theta}$, characterizing the model. By **Bayes' theorem**, we have that

$$\pi(\boldsymbol{\theta} | \mathbf{Y}, \mathbf{X}) = \frac{\mathcal{L}(\mathbf{Y} | \mathbf{X}, \boldsymbol{\theta})\pi(\boldsymbol{\theta})}{\int_{\Theta} \mathcal{L}(\mathbf{Y} | \mathbf{X}, \boldsymbol{\theta})\pi(d\boldsymbol{\theta})}$$

We can use the previous to perform models estimation and inference. Usually we don't have **conjugacy** \Rightarrow the estimation can be tedious, and computational intensive!

4/20

Mixed or mixture?

The difference

By now, we set $\boldsymbol{\theta} = (\boldsymbol{\beta}, \boldsymbol{\lambda}, \boldsymbol{\xi}) \in \Theta$.

Assume that our data are divided into **groups**. Let G a variable denoting the **allocation** of the observations to a specific group, i.e. $G_i = j \Rightarrow$ the observation i -th belongs to the group j .

- ▶ **Mixed models:** the values G_1, \dots, G_n are known and fixed, denoting a group structure of the data. As consequence, our data are **partially exchangeable** (see the next slide). GLM within this framework are defined by a relation as the following

$$g(\mu) = \mathbf{x}^\top \boldsymbol{\beta} + \mathbf{z}^\top \boldsymbol{\lambda}_j$$

where \mathbf{x} denotes the covariates for the fixed effects, $\boldsymbol{\beta}$ the coefficients for the fixed effects, \mathbf{z} denotes the covariates for the random effects, and $\boldsymbol{\lambda}$ the coefficients for the random effects.

- ▶ **Mixture models:** the values G_1, \dots, G_n are unknown and random, object of interest for our inference. The data are exchangeable.

5/20

Partial exchangeability

\rightarrow observations are exchangeable within groups

Let K be the number of groups, and $Y_{1,1}, \dots, Y_{n_1,1}, \dots, Y_{1,K}, \dots, Y_{n_K,K}$ a set of data divided into K groups.

We say that $Y_{1,1}, \dots, Y_{n_1,1}, \dots, Y_{1,K}, \dots, Y_{n_K,K}$ are **partially exchangeable** if for any set of k permutation $\sigma_1, \dots, \sigma_k$, where σ_j is a permutation of $1, \dots, n_j$, we have

$$\begin{aligned} \mathcal{L}(Y_{1,1}, \dots, Y_{n_1,1}, \dots, Y_{1,K}, \dots, Y_{n_K,K}) \\ = \mathcal{L}(Y_{\sigma_1(1),1}, \dots, Y_{\sigma_1(n_1),1}, \dots, Y_{\sigma_K(1),K}, \dots, Y_{\sigma_K(n_K),K}) \end{aligned}$$

- ▶ De Finetti representation theorem for partially exchangeable data
- ▶ The likelihood becomes

$$\begin{aligned} \mathcal{L}(Y_1, \dots, Y_K | \mathbf{X}, \mathbf{G}, \boldsymbol{\beta}, \boldsymbol{\lambda}_1, \dots, \boldsymbol{\lambda}_K) \\ = \prod_{j=1}^K \prod_{i: G_i=j}^n f_Y(Y_i | \mathbf{x}_i, \mathbf{z}_i, \boldsymbol{\beta}, \boldsymbol{\lambda}_j) \end{aligned}$$

6/20

GLMM

GLMM - introduction

! Our data are divided into K **groups**. We know the group allocation. We observe:

- ▶ $Y_{i,j}$ i -th response variable in the j -th group
- ▶ $x_{i,j}$ i -th covariates in the j -th group for the **fixed effects** → associated to global coefficients
- ▶ $z_{i,j}$ i -th covariates in the j -th group for the **random effects** → associated to coefficients which are group specific

We denote by

- ▶ $\pi(\beta)$ the prior for the regression coefficients corresponding to the **fixed effects** β
- ▶ $\prod_{j=1}^K \pi_j(\lambda_j)$ the prior for the regression coefficients corresponding to the **random effects** λ_j , for the groups $j = 1, \dots, K$
- ▶ $\pi(\xi)$ the prior for the remaining parameters of the model

7/20

GLMM - model specification

Model specification:

$$\begin{aligned} Y_{i,j} | \mu_{i,j}, \xi &\sim f_Y(Y_{i,j} | \mu_{i,j}, \xi), & i = 1 : n_j, j = 1 : K \\ \mu_{i,j} &= g^{-1}(x_{i,j}^\top \beta + z_{i,j}^\top \lambda_j) & i = 1 : n_j, j = 1 : K \\ \lambda_j &\sim \pi_j(\lambda_j), & j = 1, \dots, K \\ \beta &\sim \pi(\beta) \\ \xi &\sim \pi(\xi) \end{aligned}$$

Prior elicitation:

- ▶ We assume independent distribution with expected value equal to 0 for all the coefficients, as for example

$$\beta_\ell \sim N(0, \sigma_{\beta_\ell}^2)$$

- ▶ We may relax the model by assuming hyperprior distributions on the variances of the priors for the coefficients.

8/20

GLMM - Identifiability

Let $\mathcal{F} = \{f(\cdot, \theta) : \theta \in \Theta\}$ our statistical model, indexed by θ . We say that a **model is identifiable** if

$$f(\cdot, \theta_1) = f(\cdot, \theta_2) \implies \theta_1 = \theta_2, \quad \forall \theta_1, \theta_2 \in \Theta$$

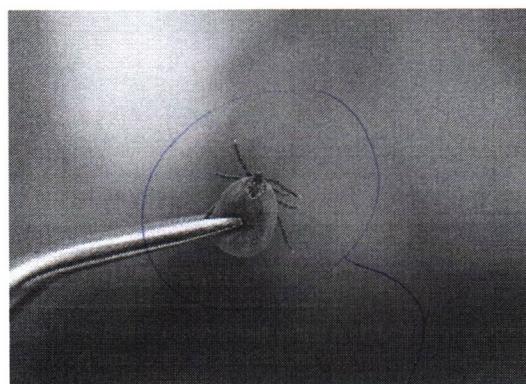
- ▶ The GLMMs are not always identifiable
- ▶ A solution to the identifiability problem is given by setting a 0-mean distribution on the random effects (as we have done)
- ▶ When you go down in the hierarchies, you can suppress the common intercept term (is helping)

9/20

Example: random intercepts

The data

A tick:



10/20

The data²

A red grouse:



11/20

The data³ and the model

We have a set of 403 observations, where each observation correspond to a particular grouse. For each animal we observe:

- ▶ Height, above sea level (meters), denoted by X , where the animal lives, continuous variable
- ▶ Year (three different years), denoted by G , where $G = 1$ is 1994, $G = 2$ is 1995 and $G = 3$ is 1996. We assume a specific random effect for each year
- ▶ Number of ticks Y on the heads of the red grouse

We assume a Poisson distribution on the data, then the **model** becomes:

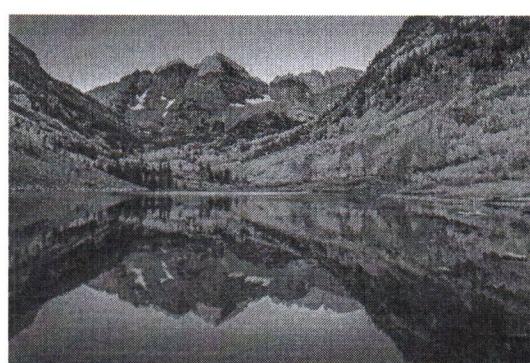
$$\begin{aligned} Y_{i,j} \mid \mu_{i,j} &\sim \text{Poisson}(\mu_{i,j}) \\ \mu_{i,j} &= e^{\beta_0 + \beta_1 x_{i,j} + \lambda_j} \\ \beta_\ell &\sim N(0, \sigma_{\beta_\ell}^2) \\ \lambda_\ell &\sim N(0, \sigma_{\lambda_\ell}^2) \\ \sigma_{\dots}^2 &\sim \text{InvGamma}(1, 10) \end{aligned}$$

12/20

Example: random coefficients

The data

Aspen:



13/20

The data²

An ant:



14/20

The data³

We have a set of 60 observations, where each observation correspond to a particular tree. The tree. are divided into High and Low altitude trees.
For each spot we observe:

- ▶ Trap days, different lengths of time surveying each tree, denoted by X
- ▶ Elevation, G . If $G = 1$, then the spot is *High*, if $G = 2$ then the spot is *Low*.
- ▶ Valley, three different valleys for each elevation, denoted by SG , nested within each elevation
- ▶ Summer precipitation (avg), assuming a random effect, denoted by Z_1
- ▶ Summer temperature (avg), assuming a random effect, denoted by Z_2
- ▶ Number of ants, assuming a random effect, denoted by Y

15/20

The model

We assume a Poisson distribution on the data, then the **model** becomes:

$$\begin{aligned}
 Y_{i,j,r} \mid \mu_{i,j,r} &\sim \text{Poisson}(\mu_{i,j,r}) \\
 \mu_{i,j,r} &= \exp(\beta_1 x_{i,j} + \theta_j + \gamma_{j,1} z_{1,i,j} + \gamma_{j,2} z_{2,i,j} + \alpha_{j,r}) \\
 \beta_1 &\sim N(0, \sigma_{\beta_1}^2) \\
 \theta_j &\sim N(0, \sigma_{\theta_j}^2) \\
 \gamma_{j,\ell} &\sim N(0, \sigma_{\gamma_{j,\ell}}^2) \\
 \alpha_{j,r} &\sim N(0, \sigma_{\alpha_{j,r}}^2) \\
 \sigma^2 &\sim \text{InvGamma}(1, 10)
 \end{aligned}$$

where we have two nested hierarchies, with β_1 common regression coefficient shared across different groups, θ_j is the random effect for the j -th group, $\alpha_{j,r}$ is the random effect of the r -th subgroup nested in the j -th group and $\gamma_{j,1}, \gamma_{j,2}$ are the random coefficients for the j -th group.

16/20

See the code!

Models for time series

Univariate time series

Let Y be a random variable, with $Y \in \mathcal{Y} \subseteq \mathbb{R}$. We further assume that we observed the random variable at different times $t = 1, \dots, T$ (for simplicity we assume discrete and equally spaced times).

Different possible modeling strategies:

- AR(p) - Auto-Regressive model of order p : = regression on the past

$$Y_t = \alpha + \sum_{j=1}^p \phi_j Y_{t-j} + \varepsilon_t, \quad (\varepsilon_t \sim f_\varepsilon) \text{ white noise}$$

- MA(p) - Moving Average model of order p : = we work with the past error terms

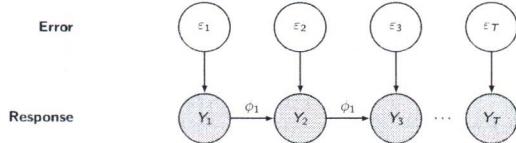
$$Y_t = \alpha + \sum_{j=1}^p \theta_j \varepsilon_{t-j}$$

- ARMA(p,q) - Auto-Regressive Moving Average model of order (p, q) , combining together the previous strategies.

17/20

AR(1)

We consider an AR(1) model.



We complete the model specification by assuming a Gaussian distribution for the error and with the following priors:

We can translate the model in a Bayesian framework:

$$Y_t | \alpha, \phi_1, \sigma_Y^2 \sim N(\alpha + \phi_1 Y_{t-1}, \sigma_Y^2)$$

$$\alpha \sim N(0, \sigma_\alpha^2)$$

$$\phi_1 \sim N(0, \sigma_\phi^2)$$

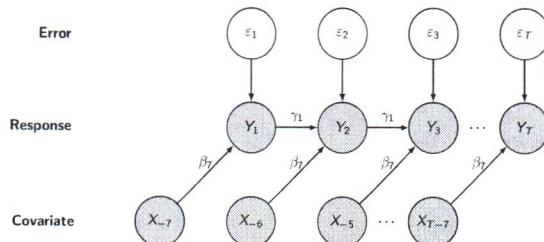
$$\sigma_Y^2 \sim InvGamma(a_Y, b_Y)$$

18/20

See the code!

DLR model

We might want to include a lagged covariate to increase the prediction performance of the model. We consider a **Dynamic Regression Model**.



19/20

DLR model

We complete the model specification assuming the following priors:

$$Y_t | \alpha, \gamma_1, \beta_1, \sigma_Y^2, \mathbf{X} \sim N(\alpha + \gamma_1 Y_{t-1} + \beta_1 X_{t-7}, \sigma_Y^2)$$

$$\alpha \sim N(0, \sigma_\alpha^2)$$

$$\phi_1 \sim N(0, \sigma_\phi^2)$$

$$\sigma_Y^2 \sim InvGamma(a_Y, b_Y)$$

20/20

See the code!

```
## S+Sp_model.stan
#####
# State-space univariate model #####
#####
# DATA ///////////////////////////////////////////////////
# data {
#   int<lower = 0> N; // number of data
#   real Y[N]; // response vector
#   real sigmaphi;
#   real phi;
#   real a_sigma2;
#   real b_sigma2;
#   real md;
# }
# MODEL ///////////////////////////////////////////////////
# model {
#   # Likelihood
#   for (t in 2:N)
#   {
#     Y[t] ~ normal(md + phi * Y[t - 1], pow(sigma2, 0.5));
#   }
#   # m0 ~ normal(0, sigma2m0);
#   phi ~ normal(0, sigma2phi);
#   sigma2 ~ inv_gamma(a_sigma2, b_sigma2);
# }
##### DLR_model.stan
#####
# State-space univariate model #####
#####
# DATA ///////////////////////////////////////////////////
# data {
#   int<lower = 0> N; // number of data
#   real Y[N]; // response vector
#   real X[N]; // response vector
#   real sigma_coeff[3];
#   real sigma2Y_par[2];
# }
# MODEL ///////////////////////////////////////////////////
# model {
#   # Likelihood
#   for (t in 2:N)
#   {
#     Y[t] ~ normal(m0 + beta * X[t] + gamma * Y[t - 1], pow(sigma2Y, 0.5));
#   }
#   # m0 ~ normal(0, sigma_coeff[1]);
#   beta ~ normal(0, sigma_coeff[2]);
#   gamma ~ normal(0, sigma_coeff[3]);
#   sigma2Y ~ inv_gamma(sigma2Y_par[1], sigma2Y_par[2]);
# }
```

Poisson.stan

```
#####
# GLM - POISSON REGRESSION #####
#####
# DATA ///////////////////////////////////////////////////
# data {
#   # int<lower = 0> N; // number of data
#   # int<lower = 0> p_fix; // number of covariates, fixed effects
#   # parameters {
#     vector[p_fix] beta; // regression coefficients
#     vector[p_fix] sigma2_beta; // variance for the prior on beta
#   }
#   # transformed parameters
#   {
#     vector[N] mu;
#     for(t in 1:N){
#       mu[t] = exp(cov(X, i) * beta);
#     }
#   }
#   # MODEL ///////////////////////////////////////////////////
# model {
#   # likelihood
#   for (s in 1:N)
#   {
#     // print(" s = ", s);
#     Y[s] ~ poisson(mu[s]);
#   }
# }
##### TRANSFORMED PARAMETERS ///////////////////////////////////////////////////
# transformed parameters
{
#   vector[N] Log_Lik;
#   for (j in 1:N)
#   {
#     Log_Lik[j] = poisson_lpmf(Y[j] | mu[j]);
#   }
# }
```

```

## Poisson_GLM.stan
# ///////////////////////////////////////////////////////////////////
# // GLMM - POISSON REGRESSION
# ///////////////////////////////////////////////////////////////////
# /////////////////////////////////////////////////////////////////// DATA ///////////////////////////////////////////////////////////////////
# data{
#   int<lower = 0> N;           // number of data
#   int<lower = 0> ngr;          // number of groups
#   int<lower = 0> nsubgr;        // number of covariates, fixed effects
#   int<lower = 0> p_fix;         // number of covariates, random effects
#   int<lower = 0> p_ran;         // response vector
#   vector[N] X;                // design matrix (fixed effects)
#   matrix[N, p_fix] Z;          // design matrix (random effects)
#   matrix[N, p_ran] Y;          // groups (dummy) allocation
# }
# /////////////////////////////////////////////////////////////////// PARAMETERS ///////////////////////////////////////////////////////////////////
# parameters{
#   vector[p_fix] beta;          // regression coefficients
#   vector[ngr] theta;            // (group specific) random effects
#   vector[p_fix] sigma2_beta;     // variances for the prior on beta
#   vector[ngr] sigma2_theta;      // variances for the prior on theta
# }
# /////////////////////////////////////////////////////////////////// TRANSFORMED PARAMETERS ///////////////////////////////////////////////////////////////////
# transformed parameters
# {
#   vector[N] mu;
#   for(i in 1:N){
#     mu[i] = exp(row(X, i) * beta + row(G, i) * theta);
#   }
# }
# /////////////////////////////////////////////////////////////////// MODEL ///////////////////////////////////////////////////////////////////
# model{
#   #// Likelihood
#   for(s in 1:N)
#   {
#     // print(" s = " , s);
#     Y[s] ~ poisson(mu[s]);
#   }
#   for(j in 1:p_fix)
#   {
#     beta[j] ~ normal(0, 0, pow(sigma2_beta[j], 0.5));
#   }
#   for(j in 1:ngr)
#   {
#     theta[j] ~ normal(0, 0, pow(sigma2_theta[j], 0.5));
#     sigma2_theta[j] ~ inv_gamma(2., 10.);
#   }
#   for(j in 1:nsubgr)
#   {
#     sigma2_beta[j] ~ normal(0, 0, pow(sigma2_beta[j], 0.5));
#   }
# }
# /////////////////////////////////////////////////////////////////// GENERATED QUANTITIES ///////////////////////////////////////////////////////////////////
# generated quantities
# {
#   vector[N] log_Lik;
#   for(j in 1:N){
#     log_Lik[j] = poisson_lpmf(Y[j] | mu[j]);
#   }
# }
# /////////////////////////////////////////////////////////////////// GENERATED QUANTITIES ///////////////////////////////////////////////////////////////////
# generated quantities{
#   vector[N] Log_Lik;
#   for(j in 1:N){
#     Log_Lik[j] = poisson_lpmf(Y[j] | mu[j]);
#   }
# }

```

```

## Poisson_GLM.stan
# ///////////////////////////////////////////////////////////////////
# // GLMM - POISSON REGRESSION
# ///////////////////////////////////////////////////////////////////
# /////////////////////////////////////////////////////////////////// DATA ///////////////////////////////////////////////////////////////////
# data{
#   int<lower = 0> N;           // number of data
#   int<lower = 0> ngr;          // number of groups
#   int<lower = 0> nsubgr;        // number of covariates, fixed effects
#   int<lower = 0> p_fix;         // number of covariates, random effects
#   int<lower = 0> p_ran;         // response vector
#   vector[N] X;                // design matrix (fixed effects)
#   matrix[N, p_fix] Z;          // design matrix (random effects)
#   int<lower = 1, upper = nsubgr> subgroup[N];
# }
# /////////////////////////////////////////////////////////////////// PARAMETERS ///////////////////////////////////////////////////////////////////
# parameters{
#   real beta;                  // regression coefficients (fixed)
#   vector[ngr] theta;            // (group specific) random effects
#   matrix[p_ran, ngr] gamma;     // regression coefficients (random)
#   vector[nsubgr] alpha;         // (group specific) random effects
#   real sigma2_beta;             // variances for the prior on beta
#   vector[ngr] sigma2_theta;      // variances for the prior on theta
#   matrix[ngr, p_ran] sigma2_gamma; // (group specific) random effects
#   vector[nsubgr] sigma2_alpha;    // (group specific) random effects
# }
# /////////////////////////////////////////////////////////////////// TRANSFORMED PARAMETERS ///////////////////////////////////////////////////////////////////
# transformed parameters{
#   vector[N] mu;
#   for(i in 1:N){
#     mu[i] = exp(X[i] * beta + theta[group[i]] + alpha[subgroup[i]]);
#   }
# }
# /////////////////////////////////////////////////////////////////// MODEL ///////////////////////////////////////////////////////////////////
# model{
#   #// Likelihood
#   for(s in 1:N){
#     Y[s] ~ poisson(mu[s]);
#   }
#   #// priors
#   beta ~ normal(0, 0, pow(sigma2_beta, 0.5));
#   sigma2_beta ~ inv_gamma(2., 10.);
#   for(j in 1:ngr){
#     theta[j] ~ normal(0, 0, pow(sigma2_theta[j], 0.5));
#     sigma2_theta[j] ~ inv_gamma(2., 10.);
#   }
#   for(j in 1:nsubgr){
#     alpha[j] ~ normal(0, 0, pow(sigma2_alpha[j], 0.5));
#     sigma2_alpha[j] ~ inv_gamma(2., 10.);
#   }
# }
# /////////////////////////////////////////////////////////////////// GENERATED QUANTITIES ///////////////////////////////////////////////////////////////////
# generated quantities{
#   vector[N] Log_Lik;
#   for(j in 1:N){
#     Log_Lik[j] = normal(0, 0, pow(sigma2_gamma[j, i], 0.5));
#     sigma2_gamma[j, i] ~ normal(0, 0, inv_gamma(2., 10.));
#   }
# }
# /////////////////////////////////////////////////////////////////// GENERATED QUANTITIES ///////////////////////////////////////////////////////////////////
# generated quantities{
#   vector[N] Log_Lik;
#   for(j in 1:N){
#     Log_Lik[j] = poisson_lpmf(Y[j] | mu[j]);
#   }
# }

```

19/11/2020

6/8

```

#####
##### GLMM #####
#####

# R & STAN are friends!
library(rstan)
library(coda)

# for plots
library(ggplot2)
library(tidyverse)
library(dplyr)
library(purrr)
library(bigrss)
require(BIGplots)
require(eggpubr)

load("data_grouse")
head(data_grouse)

X <- data_grouse[,2]
temp <- as.factor(data_grouse$YEAR)
G <- model.matrix(~temp, contrasts=list(temp=contrasts(temp, contrasts=F)))[,-1]
colnames(G) <- c("95", "96", "97")
Y <- as.vector(data_grouse[,1])

N <- length(X)
p_fix <- 2
ng <- length(unique(data_grouse$YEAR))

# useful function
my_WAIC <- function(fit, param){
  rstan::extract(fit, param)[[1]]
}

## WAIC <- sum(sapply(WAIC, 2, var))
## pdpd <- sum(sapply(WAIC, 2, function(x) log(mean(exp(x)))))
## WAIC <- -2 * pdpd + 2 * p_WAIC
## return(WAIC)
}

#####
##### POI regression #####
#####

data_POI <- list(N = N,
                  p_fix = 2,
                  ng = 3,
                  Y = Y,
                  X = as.matrix(cbind(rep(1, N), X)),
                  G = as.matrix(G))

inits <- function()
{
  list(beta = rep(0,1, p_fix),
       data = data_POI,
       chains = 2,
       iter = 30000,
       warmup = 10000,
       thin = 20,
       seed = 42,
       init = inits,
       algorithm = 'NUTS')
}

# save(Poi_GLMM, file="Poi_GLMM_est.dat")
load("Poi_GLMM_est.dat")

rstan::traceplot(Poi_GLMM, pars = "sigma2_theta", inc_warmup = TRUE)

```

```

ylab("Number of ticks on red grouse") +
xlab("")
```

```

#----#
## COMPARE WITH GLM
#----#
data_POI <- list(N = N,
                  p_fix = 2,
                  Y = Y,
                  X = as.matrix(cbind(rep(1, N), X)))
```

```

inits <- function()
{
  list(beta = rep(0,1, p_fix),
       sigma2_beta = rep(10, p_fix))
}
```

```

##### FIT THE MODEL
```

```

Poi_GLM <- stan(file = "Poi_GLM.stan",
                 data = data_POI,
                 chains = 2,
                 iter = 300000,
                 warmup = 100000,
                 thin= 20,
                 seed = 42,
                 init = inits,
                 algorithm = 'NUTS')
save(Poi_GLM, file="Poi_GLM.est.dat")
load("Poi_GLM.est.dat")
```

```

rstan::traceplot(Poi_GLM, pars = "sigma2_beta", inc_warmup = TRUE)
rstan::traceplot(Poi_GLM, pars = "beta", inc_warmup = FALSE)
```

```

# Diagnostic -----
```

```

coda_chain <- As.mcmc.list(Poi_GLM, pars = c("beta"))
summary(coda_chain)
```

```

# Compute the gelman and rubin's convergence diagnostic
# ADB: potential scale reduction factor
gelman.diag(coda_chain, confidence = 0.95, autoburnin = TRUE, multivariate=TRUE)
```

```

# Geweke's convergence diagnostic
geweke.diag(coda_chain, frac1=0.1, frac2=0.5)
```

```

# autocorrelation and plot
acfplot(coda_chain, lag.max = 30)
```

```

# WAIC GLMM VS GLM -----
my_WAIC(Poi_GLM, "log_lik")
my_WAIC(Poi_GLM, "log_lik")
```

```

# -----
```

```

# -----
```

```

# -----
```

```

AntData <- read.csv("AntAbundance.csv")
head(AntData)
```

```

# num.fpods ~ b0..j * elevation + b1..j * valley + b2..j * summer.precip + b3..j * summer.temp + b4 * trap.da
ys + table(AntData$num.fpods)
```

```

Y <- AntData$num.fpods
X <- as.matrix(AntData$trap.days)
group <- as.numeric(as.factor(AntData$elevation))
subgroup <- rep(1:6, each = 18)
Z <- as.matrix(cbind(AntData[, c(3,10)]))
N <- length(Y)
```

```

p_fix <- 1
```

```

plot_post <- Poi_GLM2$stan$parameters
rstan::extract("phi") %>%
  map_df(as.data.frame, .id = 'param') +
  theme_minimal() +
  theme(legend.position="none")

# theta
plot_post %>%
  rstan::extract("theta") %>%
  map_df(as.data.frame, .id = 'param') +
  geom_density() +
  facet_wrap(~param, scales = 'free') +
  scale_fill_locuszoom() +
  theme_minimal() +
  theme(legend.position="none")

plot_post %>%
  rstan::extract("gamma") %>%
  map_df(as.data.frame, .id = 'param') +
  geom_density() +
  facet_wrap(~param, scales = 'free') +
  scale_fill_locuszoom() +
  theme_minimal() +
  theme(legend.position="none")

plot_post %>%
  rstan::extract("lambda") %>%
  map_df(as.data.frame, .id = 'param') +
  geom_density() +
  facet_wrap(~param, scales = 'free') +
  scale_fill_locuszoom() +
  theme_minimal() +
  cord_flip() +
  xlim(c(-180, 20)) +
  theme(legend.position="none")

# alpha
plot_post <- Poi_GLM2$stan$parameters
rstan::extract("alpha") %>%
  as.data.frame() %>%
  map_df(as.data.frame, .id = 'param') +
  theme_minimal() +
  theme(legend.position="none")

plot_post %>%
  rstan::extract("beta") %>%
  map_df(as.data.frame, .id = 'param') +
  geom_density() +
  facet_wrap(~param, scales = 'free') +
  scale_fill_locuszoom() +
  theme_minimal() +
  cord_flip() +
  xlim(c(-10,8)) +
  theme(legend.position="none")

# theta
plot_post <- Poi_GLM2$stan$parameters
rstan::extract("phi") %>%
  as.data.frame() %>%
  map_df(as.data.frame, .id = 'param') +
  geom_density() +
  facet_wrap(~param, scales = 'free') +
  scale_fill_locuszoom() +
  theme_minimal() +
  theme(legend.position="none")

plot_post %>%
  rstan::extract("gamma") %>%
  map_df(as.data.frame, .id = 'param') +
  geom_density() +
  facet_wrap(~param, scales = 'free') +
  scale_fill_locuszoom() +
  theme_minimal() +
  theme(legend.position="none")

# alpha
plot_post <- Poi_GLM2$stan$parameters
rstan::extract("alpha") %>%
  as.data.frame() %>%
  map_df(as.data.frame, .id = 'param') +
  geom_density() +
  facet_wrap(~param, scales = 'free') +
  scale_fill_locuszoom() +
  theme_minimal() +
  cord_flip() +
  xlim(c(-1, 12)) +
  theme(legend.position="none")

# beta
plot_post <- Poi_GLM2$stan$parameters
rstan::extract("beta") %>%
  as.data.frame() %>%
  map_df(as.data.frame, .id = 'param') +
  geom_density() +
  facet_wrap(~param, scales = 'free') +
  scale_fill_locuszoom() +
  theme_minimal() +
  theme(legend.position="none")

# Compare with a GLM (Poisson regression - no valley)
AntData$POI <- list(N = length(Y),
                     D_fix = p,
                     Y = Y,
                     X = as.matrix(X))
AntData$POI

# Compare with a GLM (Poisson regression - no valley)
AntData$POI <- list(N = length(Y),
                     D_fix = p,
                     Y = Y,
                     X = as.matrix(X))
AntData$POI

# Fit the model
StSp_model <- stan(file = "StSp_model.stan",
                     data = data_StSp,
                     chains = 2,
                     iter = 30000,
                     warmup = 10000,
                     thin= 20,
                     seed = 42,
                     init = 1)

# Fit the model
StSp_model <- stan(file = "StSp_model.est.stan",
                     data = data_StSp,
                     chains = 2,
                     iter = 30000,
                     warmup = 10000,
                     thin= 20,
                     seed = 42,
                     init = 1)

# Fit the model
StSp_model <- stan(file = "StSp_model.est.dat")
load("StSp_model.est.dat")
load("StSp_model.est.dat")

# Fit the model
StSp_model <- stan(file = "StSp_model",
                     data = data_StSp,
                     algorithm = 'NUTS')

```



```

} else {
  resultDLR[,1] <- apply(cbind(params, resultLR[,i-1]), 1, function(x) rnorm(1, x[2] + x[1] * Xtemp
  [i] + x[3] * x[5], sqrt(x[4])))
}

ggplot(data.frame(x = (length(Y)):(length(Y) + 6),
  y = colMeans(resultDLR),
  y_low = apply(resultDLR, 2, quantile, p = 0.05),
  y_high = apply(resultDLR, 2, quantile, p = 0.95))) +
  geom_line(data = data.frame(x = (length(Y)):(length(Y) + 6),
    y = colMeans(resultSSp),
    y_low = apply(resultSSp, 2, quantile, p = 0.05),
    y_high = apply(resultSSp, 2, quantile, p = 0.95)),
  aes(x = X, y = Y), lty = 2, col = 3) +
  geom_ribbon(data = data.frame(x = (length(Y)):(length(Y) + 6),
    y = colMeans(resultSSp),
    y_low = apply(resultSSp, 2, quantile, p = 0.05),
    y_high = apply(resultSSp, 2, quantile, p = 0.95)),
  aes(ymin = y_low, ymax = y_high, x = X), alpha = 0.2, fill = 3) +
  geom_line(data = data.frame(x = X, y = Y), lty = 2) + theme_bw() +
  geom_ribbon(aes(x = X, y = Y), x = 1:length(Y), aes(x = X, y = Y)) +
  geom_vline(aes(xintercept = length(Y), col = 2, lty = 2) +
  ggtitle("Aki model & DLR - COVID data Lombardia"))

```

Survival data

Riccardo Corradin
November 27, 2020

1/23

Survival data

Introduction

Let T be a **RV** denoting a positive time to event, i.e. with support \mathbb{R}^+ , object of interest. Further we denote by $F(t) = P[T \leq t]$ the **cdf** of T and by $f(t) = \frac{d}{dt}F(t)$ its **pdf**.

Related quantities of interest

- ▶ Survival function - surviving at least until t

$$S(t) = P[T > t] = 1 - F(t)$$

- ▶ Hazard function - event rate at t conditionally on survival until t

$$h(t) = \lim_{\delta t \rightarrow 0} \frac{P[t \leq T < t + \delta t \mid T \geq t]}{\delta t}$$

= probability that the event is happening during this infinitesimal t (δt)

- ▶ Cumulative hazard function - cumulative hazard until time t

$$H(t) = \int_0^t h(x)dx$$

2/23

Introduction²

Useful relations

- ▶ Hazard - survival

$$h(t) = -\frac{S'(t)}{S(t)} = -\frac{d}{dt} \log(S(t))$$

- ▶ Survival - cumulative hazard

$$S(t) = e^{-H(t)}$$

- ▶ Density - survival & hazard

$$f(t) = h(t)S(t)$$

Goal: we want to model the time to event, i.e. the survival function related to T or the cdf.

3/23

Right censoring

In an perfect world, we run the studies until the end of "Life, the Universe and Everything" and we observe the occurrences of all the events. But in real world, we do not.

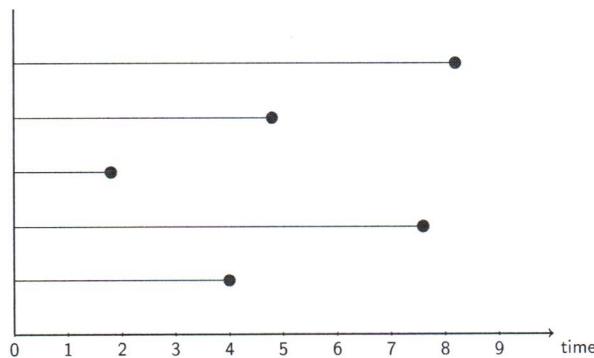
In practice, three scenarios may occur:

- An event occurs
- We are monitoring an individual/thing/what-ever, but the event does not occur before the end of the study
- We are monitoring an individual/thing/what-ever, but is leaving the study before the event occurs

Suppose that we observe T_1, \dots, T_n of time to events.

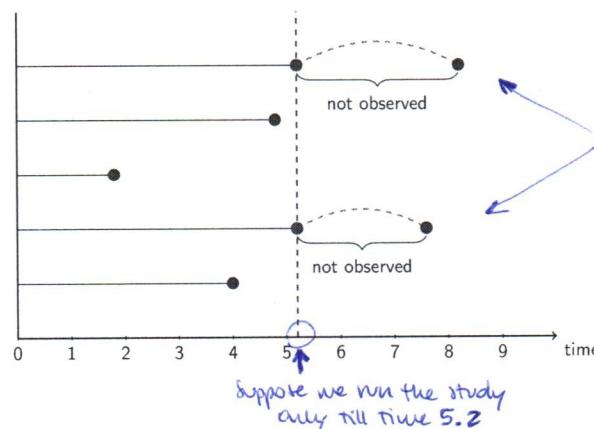
4/23

Right censoring²



5/23

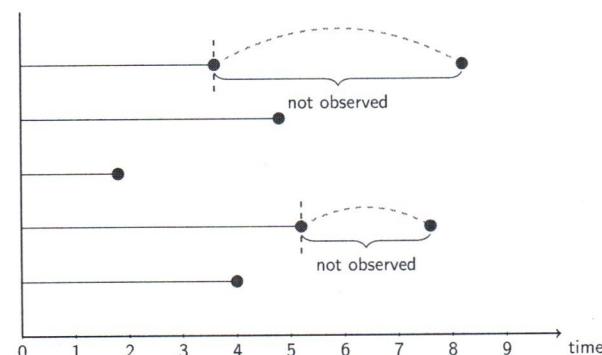
Right censoring³ : FIXED CENSORING



6/23

Right censoring⁴

: RANDOM CENSORING = the censoring time is different individual to individual



7/23

The likelihood

For simplicity we assume a common censoring time c , such that $T_i = \min\{Y_i, c\}$. Let δ_i denotes if an observation is censored, i.e.

$$\delta_i = \begin{cases} 1 & Y_i \leq c \\ 0 & Y_i > c \end{cases}$$

We denote by the generic $g(T_i | \cdot)$ the contribution to the likelihood of the i -th observation

- Each observation not truncated contributes to the likelihood by the fact that it survived until T_i and the event occurs at T_i , i.e.

$$g(T_i | \delta_i = 1, \theta) = f(T_i, \theta) = h(T_i, \theta)S(T_i, \theta)$$

- Each truncated observation contributes to the likelihood by the fact that the event occurs after T_i (i.e. survived until T_i without occurrence)

$$g(T_i | \delta_i = 0, \theta) = S(T_i, \theta)$$

8/23



The likelihood²

Within the previous setup, the likelihood of an n -sized sample becomes

$$\begin{aligned} \mathcal{L}(T_1, \dots, T_n | \delta_1, \dots, \delta_n, \theta) &= \prod_{i=1}^n g(T_i | \delta_i, \theta) \\ &= \prod_{i=1}^n [f(T_i | \theta)]^{\delta_i} [S(T_i | \theta)]^{1-\delta_i} \\ &= \prod_{i=1}^n S(T_i | \theta)h(T_i | \theta)^{\delta_i} \end{aligned}$$

where it is apparent the contribution of the two situations:

- all the observations contributed in terms of surviving until the last observed time
- the uncensored observations contribute also by the risk of observe the events

9/23

Different parametric models

Different choices for the parametric model $f(T | \theta)$ lead to different models for the time to events. Common choices:

- **Exponential** $f(T_i | \theta) = \theta e^{-\theta T_i} \quad T_i > 0, \theta > 0$
- **Weibull** $f(T_i | \theta, \kappa) = \theta \kappa t^{\kappa-1} e^{-\theta t^\kappa} \quad T_i > 0, \theta > 0, \kappa > 0$
- **Log-logistic**

$$f(T_i | \alpha, \beta) = \frac{\beta T_i^{\beta-1} / \alpha^\beta}{(1 + (T_i/\alpha)^\beta)^2} \quad T_i > 0, \alpha > 0, \beta > 0$$
- **Log-normal** $f(T_i | \mu, \sigma) = \frac{1}{T_i \sigma \sqrt{2\pi}} \exp\left(-\frac{1}{2} \left(\frac{\log T_i - \mu}{\sigma}\right)^2\right) \quad T_i > 0, \mu \in \mathbb{R}, \sigma > 0$

10/23

Example: weights

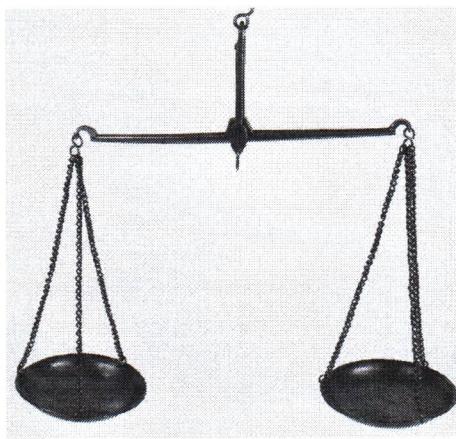
— here we decide to treat the censored date as missing values date (instead of doing *)

A potato



11/23

A scale



12/23

The weights example

It is not only about the time, we can censor other stuffs! For example, the same mechanism works also if you consider the weights problem:

- ▶ We have a set of potatoes, with weights distributed as $W_i \sim \text{Exp}(\theta)$
- ▶ We have a scale, with an upper limit c
- ▶ We have censored observations for the potatoes with weights larger than c

We are interested about inference for the parameter θ

- ▶ **Conjugate prior** $\pi(\theta) \sim \text{Gamma}(a_0, b_0)$
- ▶ **Likelihood**

$$\mathcal{L}(W_1, \dots, W_n | \delta_1, \dots, \delta_n, \theta) = \prod_{i: \delta_i=1} \theta e^{-\theta W_i} \times \prod_{i: \delta_i=0} e^{-\theta W_i}$$

13/23

Augment the problem

We have a sample $\{(W_1, \delta_1), \dots, (W_n, \delta_n)\}$ of censored weights and censor variables.

- ▶ If the i -th observation is censored, we do not know what is the i -th weight
- ▶ We handle the i -th weight as a missing value

Core idea behind the algorithm:
a **Gibbs sampler** for $\pi(\theta, \mathbf{W}^{(miss)} | \mathbf{W}^{(obs)})$ where

1. We sample $\mathbf{W}^{(miss)} \sim f(W | \theta)$
2. We sample $\theta \sim \pi(\theta | \mathbf{W}^{(miss)}, \mathbf{W}^{(obs)})$

14/23

See the code!

Example 2. Kevlar fiber



15/23

Spool



16/23

The data

We have a set of data composed by 108 kevlar fiber lifetimes, coming from the combination of eight different spools and four levels of stress (pressure).

A subset of the data (11 observations) are censored.

We are interested into model the lifetimes of the kevlar fibers, taking into account the different spools and the stress levels.

First let's see some preliminary analysis..

See the code!

PH vs AFT models

Two main family of models for survival data

- ▶ **Proportional hazard models:** the effect of the covariates is multiplicative on the hazard rate

$$h(T_i | \beta, \mathbf{x}_i) = e^{\beta^T \mathbf{x}_i} h_0(t_i)$$

where $h_0(t_i)$ is a baseline hazard rate

- ▶ **Accelerated failure time models:** the effect of the covariate is multiplicative in the lifetime, which in terms of the hazard rate can be written as

$$h(T_i | \beta, \mathbf{x}_i) = e^{\beta^T \mathbf{x}_i} h_0(e^{\beta^T \mathbf{x}_i} t_i)$$

If $\beta^T \mathbf{x}_i < 0$ the effect of the covariates is to decelerate the time, whereas a $\beta^T \mathbf{x}_i > 0$ accelerates the time

} distortion at hazard level

} distortion at hazard level and lifetime level (time level)

18/23

The Weibull case

We consider a **Weibull distribution** with the following parametrization

$$X \sim \text{Wei}(\kappa, \theta), \quad \kappa, \theta \in \mathbb{R}^+$$

and the following related quantities

$$\begin{aligned} f(x; \kappa, \theta) &= \frac{\kappa}{\theta} \left(\frac{x}{\theta} \right)^{\kappa-1} e^{-(\frac{x}{\theta})^\kappa} \\ F(x; \kappa, \theta) &= 1 - e^{-(\frac{x}{\theta})^\kappa} \\ S(x; \kappa, \theta) &= e^{-(\frac{x}{\theta})^\kappa} \\ h(x; \kappa, \theta) &= \frac{\kappa}{\theta} \left(\frac{x}{\theta} \right)^{\kappa-1} \\ \mathbb{E}[X] &= \theta \Gamma \left(1 + \frac{1}{\kappa} \right) \end{aligned}$$

19/23

The Weibull case

When we use a **Weibull baseline hazard** $h_0(t)$, i.e.

$$h_0(t) = \frac{\kappa}{\theta^\kappa} t^{\kappa-1}$$

} in this case the two approaches (proportional hazard models and accelerated failure time models) coincide !

both the proportional hazards and the accelerated life models coincide, provided the covariates are constant (see Cox and Oakes, 1984, p.71).

Using for example a **Proportional hazard models** we have

$$h(T_i | \beta, \mathbf{x}_i) = e^{\beta^T \mathbf{x}_i} \frac{\kappa}{\theta^\kappa} t^{\kappa-1} = e^{\beta^T \mathbf{x}_i - \kappa \log(\theta)} \kappa t^{\kappa-1}$$

A multiplicative effect on the hazard rate is equivalent to an effect on the parameter θ as in the previous, where now we have that

$$T_i \sim \text{Weibull}(\kappa, e^{\mu - \frac{\beta^T \mathbf{x}_i}{\kappa}})$$

by setting $\mu = \log(\theta)$, with $\mu \in \mathbb{R}$.

20/23

The Weibull case: the model

We have data for different spools, we use a random effects model, with a specific intercept for each spool. **Model specification :**

$$\begin{aligned} T_i | g_j, x_i, \beta, \mu_j, \kappa &\sim \text{Weibull} \left(\kappa, e^{\mu_j - \frac{\beta^T x_i}{\kappa}} \right) & i = 1, \dots, n_j, j = 1, \dots G \\ \mu_j | m_0, \sigma_\mu^2 &\sim N(m_0, \sigma_\mu^2), & j = 1, \dots, G \\ \beta | b_0, \sigma_\beta^2 &\sim N(b_0, \sigma_\beta^2) \\ \sigma_\beta^2 | a_\beta, b_\beta &\sim IG(a_\beta, b_\beta) \\ \kappa | a_\kappa, b_\kappa &\sim Gamma(a_\kappa, b_\kappa) \end{aligned}$$

We compare the model with

- the case without random effects
- the case without covariates (only intercept)

21/23

See the code!

Example 3. The data

We have a set of data composed by 456 people, enrolled in drug addiction rehabilitation treatment.

A subset of the data (111 observations) are censored.

We are interested into model the time to events, represented by the possibly relapse of the drug. People lost to follow up are considered as events, people out of the study without events as censored.

We know the drugs usage history of the patients, and we have four different groups: Heroin and Cocaine, Heroin only, Cocaine only, no one of the previous. We use these groups to include random effects in the model.

We want to include in the model also the age and the length of treatment as covariates, possibly with random regression coefficients.

22/23

The Weibull case: the model

We use a random effects model with random regression coefficients. **Model specification :**

$$\begin{aligned} T_i | g_j, x_i, \beta, \mu_j, \kappa &\sim \text{Weibull} \left(\kappa, e^{\mu_j - \frac{\beta^T x_i}{\kappa}} \right) & i = 1, \dots, n_j, j = 1, \dots G \\ \mu_j | m_0, \sigma_\mu^2 &\sim N(m_0, \sigma_\mu^2), & j = 1, \dots, G \\ \beta_j | b_0, \sigma_\beta^2 &\sim N(b_0, \sigma_\beta^2), & j = 1, \dots, G \\ \sigma_\beta^2 | a_\beta, b_\beta &\sim IG(a_\beta, b_\beta) \\ \kappa | a_\kappa, b_\kappa &\sim Gamma(a_\kappa, b_\kappa) \end{aligned}$$

We compare the model with

- the case without random effects
- the case without covariates (only intercept)

23/23

Drugs_fixed.stan

```

data {
    int<lower = 0> N; // number of observations
    int<lower = 0> p; // number of covariates
    int<lower = 0, upper = 1> cens[N]; // censoring
    real<lower = 0> time[N]; // Failure times
    matrix[N, p] X; // design matrix
}

parameters {
    vector[p] beta; // coefficient for log(stress)
    real mu; // random effects
    real<lower=0> lambda; // variance on the mu
    real<lower=0> kappa; // second parameter of the Weibull (shape)
}

model {
    // likelihood failure times as exact or censored
    for ( j in 1:N )
    {
        if (cens[j] == 0)
        {
            time[j] ~ weibull(kappa, exp(mu));
        }
        else
        {
            target += weibull_lccdf(time[j] | kappa, exp(mu));
        }
    }

    mu ~ normal(0, pow(lambda, 0.5)); // iid
    lambda ~ inv_gamma(5, 10);
    kappa ~ gamma(1, 0.2);
}

generated quantities {
    real log_lik[N];
    for (i in 1:N)
    {
        if (cens[i] == 0)
        {
            log_lik[i] = weibull_lpdf(time[i] | kappa, exp(mu));
        }
        else
        {
            log_lik[i] = weibull_lccdf(time[i] | kappa, exp(mu));
        }
    }
}

```

Hide

Hide

Hide

```

## kevlar_dummy.stan

data {
  int<lower = 0> N; // number of observations
  int<lower = 0> ngr; // number of groups
  int<lower = 0> p; // number of covariates
  int<lower = 0, upper = 1> cens[N]; // censoring
  real<lower = 0> time[N]; // real times
  real<lower = 0> mu; // variance on the mu
  real<lower = 0> lambda; // second parameter of the Weibull (shape)
}

parameters {
  real mu; // random effects
  real<lower=0> lambda; // variance on the mu
  real<lower=0> kappa; // second parameter of the Weibull (shape)
}

model {
  // likelihood failure times as exact or censored
  for ( j in 1:N ) {
    if (cens[j] == 0) {
      time[j] ~ weibull(kappa, exp(mu));
    } else {
      target += weibull_lccdf(time[j] | kappa, exp(mu));
    }
  }
}

model {
  // coefficient for log(stress)
  vector[ngr] mu; // random effects for stress
  real<lower=0> lambda; // variance on the mu
  real<lower=0> kappa; // second parameter of the Weibull (shape)
}

parameters {
  matrix[p, ngr] beta; // coefficient for log(stress)
}

model {
  // likelihood failure times as exact or censored
  for ( j in 1:N )
  {
    if (cens[j] == 0)
    {
      time[j] ~ weibull(kappa, exp(mu[group[j]] - row(X, j) * col(beta, group[j]) / kappa));
    }
    else
    {
      target += weibull_lccdf(time[j] | kappa, exp(mu[group[j]] - row(X, j) * col(beta, group[j]) / kappa));
    }
  }
}

generated quantities {
  real log_lik[N];
  for (i in 1:N) {
    if (cens[i] == 0) {
      log_lik[i] = weibull_lpdf(time[i] | kappa, exp(mu));
    } else {
      log_lik[i] = weibull_lccdf(time[i] | kappa, exp(mu));
    }
  }
}

generated quantities {
  real log_lik[N];
  for (i in 1:N) {
    if (cens[i] == 0) {
      log_lik[i] = weibull_lpdf(time[i] | kappa, exp(mu[group[i]] - row(X, i) * col(beta, group[i])));
    } else {
      log_lik[i] = weibull_lccdf(time[i] | kappa, exp(mu[group[i]] - row(X, i) * col(beta, group[i])));
    }
  }
}

```

kevlar_random.stan

```

## kevlar_random.stan

data {
  int<lower = 0> N; // number of observations
  int<lower = 0, upper = 1> cens[N]; // censoring
  real<lower = 0> time[N]; // Failure times
  real<lower = 0> stress[N]; // Pressure
}

parameters {
  real bet; // coefficient for log(stress)
  real mu; // fixed effects
  real<lower=0> lambda; // variance on the mu
  real<lower=0> kappa; // second parameter of the Weibull (shape)
}

model {
  // likelihood failure times as exact or censored
  for ( j in 1:N )
  {
    if (cens[j] == 0)
    {
      time[j] ~ weibull(kappa, exp(mu - beta*log(stress[j]) / kappa));
    }
    else
    {
      target += weibull_lccdf(time[j] | kappa, exp(mu - beta*log(stress[j]) / kappa));
    }
  }
  // right-censored data
  target += weibull_lccdf(time[N] | kappa, exp(mu - beta*log(stress[N]) / kappa)); // right-cen
}

generated quantities {
  real log_lik[N];
  for ( i in 1:N ) {
    if (cens[i] == 0)
    {
      log_lik[i] = weibull_lpdf(time[i] | kappa, exp(mu - beta*log(stress[i]) / kappa));
    }
    else
    {
      log_lik[i] = weibull_lccdf(time[i] | kappa, exp(mu - beta*log(stress[i]) / kappa));
    }
  }
}

## kevlar_fixed.stan

data {
  int<lower = 0> N; // number of observations
  int<lower = 0, upper = 1> cens[N]; // censoring
  real<lower = 0> time[N]; // Failure times
  int<lower = 1, upper = 8> spool[N]; // Spool indicators
  real<lower = 0> stress[N]; // Pressure
}

parameters {
  real bet; // coefficient for log(stress)
  real mu; // random effects
  real<lower=0> lambda; // variance on the mu
  real<lower=0> kappa; // second parameter of the Weibull (shape)
}

model {
  // likelihood failure times as exact or censored
  for ( j in 1:N )
  {
    if (cens[j] == 0)
    {
      time[j] ~ weibull(kappa, exp(mu[ spool[j]] - beta*log(stress[j]) / kappa));
    }
    else
    {
      target += weibull_lccdf(time[j] | kappa, exp(mu[ spool[j]] - beta*log(stress[j]) / kappa));
    }
  }
  // right-censored data
  target += weibull_lccdf(time[N] | kappa, exp(mu[ spool[N]] - beta*log(stress[N]) / kappa));
}

generated quantities {
  real log_lik[N];
  for ( i in 1:N ) {
    if (cens[i] == 0)
    {
      log_lik[i] = weibull_lpdf(time[i] | kappa, exp(mu[ spool[i]] - beta*log(stress[i]) / kappa));
    }
  }
}

```

Hide

```

# R & STAN are friends!
library(rstan)
library(coda)
library(survival)

# for plots
library(ggplot2)
library(tidyverse)
library(dplyr)
library(purrr)
library(ggpubr)
library(ggplots)
require(ggpubr)

# useful function
my_WAIC <- function(fit, param){
  llik     <- rstan::extract(fit, param)[[1]]
  p_WAIC <- sum(sapply(llik, 2, var))
  WAIC   <- -2 * lppd + 2 * p_WAIC
  return(WAIC)
}

# example on censored data:
load("potatoes.dat")
head(potatoes)

gibbs_surv <- function(data, cens, niter, burnin, thin, a0, b0){

  # initialize the quantities
  results <- c()
  an <- a0 + length(data)
  bn <- b0 + sum(data)

  pb <- txtProgressBar(0, niter, style = 3)
  for(i in 1:niter){
    # sample theta from the posterior
    theta <- rgamma(1, an, rate = bn)

    # save the results
    if(i > burnin & ((i - burnin) %% thin == 0)){
      results[(i - burnin)/thin] <- theta
    }
    setTxtProgressBar(pb, i)
  }
  close(pb)
  return(results)
}

# save the results
if(i > burnin & ((i - burnin) %% thin == 0)){
  results[(i - burnin)/thin] <- theta
  data_cens[(i - burnin)/thin, ] <- aug_data[cens]
}

# sample the missing data
aug_data[cens] <- data[cens] + rexp(sum(!cens), rate = theta)

# SAMPLE A CHAIN
# ab = "prior sample size"
# bb = "prior guess on theta"
sample1 <- gibbs_surv(potatoes$potatoes, potatoes$cens, 10000, 1000, 5, 1, 50)
sample_coda <- as.mcmc(sample1[[1]])

# DIAGNOSTIC
summary(sample_coda)
plot(sample_coda)
acfplot(sample_coda, frac1=0.1, frac2=0.5)

# POSTERIOR
p2 <- ggplot(data.frame(value = sample2), aes(value)) +
  geom_density(fill = "blue", alpha = 0.5) +
  theme_minimal() +
  theme(legend.position="none") +
  geom_vline(aes(xintercept = 1 / mean(potatoes$potatoes)), col = 2, lty = 2)

# DIAGNOSTIC
summary(sample_coda)
plot(sample_coda)

```

```

summary(fit12)

# -----
# plot the residuals :
# no Bayes :
# -----
dati = read.table("kevlar.txt", header = T)
dati$atm1 = log(dati$time / 10^3) - fit2$linear
w = fit2$scale * (log(dati$time / 10^3) - fit2$linear)
par(mfrow = c(1,2))
plot(w)
abline(h = 0, col = "red")
title(paste("Generalized residuals"))
hist(w, nclass = "fd")
# -----
# -----
# KEVALAR DATA
# with Bayes :
# -----
# -----
# dev.off()

# -----
# Divide Time and stress for a constant to avoid numerical problems
dati_list = list(N = 108, time = Time(10^3),
spool = Spool,
cens = Crep(97, rep(1,11)),
stress = Stress/(23.4))

hist(log(Time/(10^3)), nclass = 30, freq = F, xlab = "time", main = "")

# -----
# STAN model
# -----
initis <- function() {
  list(beta = 1, kappa = 2, lambda = 10, mu = rep(0,8))
}

model = stan(file = "kevlar-random.stan",
data = dati_list,
chains = 1,
iter = 50000,
warmup = 10000,
thin=10,
init = initis,
control = list(max_treedepth = 20, adapt_delta = 0.99))

# save(model, file="model.dat")
load("model.dat")

# -----
# Final iterations: 4000
# Posterior mean:
get_posterior_mean(model, pars = c("beta", "mu", "lambda", "kappa"))

# -----
# What about the chains
rstan:::traceplot(model, pars = "beta", inc_warmup = TRUE)
rstan:::traceplot(model, pars = "mu", inc_warmup = TRUE)
rstan:::traceplot(model, pars = "lambda", inc_warmup = TRUE)
rstan:::traceplot(model, pars = "kappa", inc_warmup = TRUE)

# -----
# Random effects
plot(model, ask = T, pars = c("mu"))

# -----
# Diagnostic
coda_chain <- As.mcmc.list(model, pars = c("beta", "mu", "lambda", "kappa"))
Summary(coda_chain)
acfplot(coda_chain)
geweke.diag(coda_chain, frac1=0.1, frac2=0.5)

# -----
# Posterior
plot_post <- model %>%
rstan:::extract(c("kappa", "lambda", "beta")) %>%
as.data.frame() %>%
map_df(as_data_frame, .id = 'param')

plot_post %>%
ggplot(aes(value, fill = param)) +
  geom_bar(stat = "identity") +
  theme_minimal()

```



```

inits <- function() {
  list(beta = 1, kappa = 2, lambda = 10)
}

model_fix = stan(file = "kevlar_fix.stan",
  data = dati_list,
  chains = 1,
  iter = 50000,
  warmup = 10000,
  thin= 10,
  seed = 1919,
  init = inits,
  control = list(max_treedepth = 20, adapt_delta = 0.99))
  save(model_fix, file="model_fix.dat")
  load("model_fix.dat")

#-----#
# DIAGNOSTIC
rstan::traceplot(model_fix, pars = c("kappa", "mu", "beta"), inc_warmup = TRUE)

#-----#
# CODA
coda_chain_fix <- As.mcmc.list(model_fix, pars = c("lambda", "mu", "kappa", "beta"))
summary(coda_chain_fix)
acfplot(coda_chain_fix)
geweke.diag(coda_chain_fix, frac1=0.1, frac2=0.5)

#-----#
# Posterior
plot_post <- model_fix %>%
  rstan::extract(c("kappa", "mu", "beta")) %>%
  acfplot(coda_chain_fix)
geweke.diag(coda_chain_fix, frac1=0.1, frac2=0.5)

#-----#
# MODEL COMPARISON - dummy
plot_post %>%
  ggplot(aes(value, fill = param)) +
  geom_density() +
  facet_wrap(~param, scales = 'free') +
  scale_fill_locuszoom() +
  theme_minimal() +
  theme(legend.position="none")

#-----#
# MODEL COMPARISON - WAIC
my_WAIC(modele, "log_lik")
my_WAIC(modele_fix, "log_lik")
my_WAIC(modele_dummy, "log_lik")

#-----#
# STAN model
inits <- function() {
  list(kappa = 2, lambda = 10)
}

model_dummy = stan(file = "kevlar_dummy.stan",
  data = dati_list,
  chains = 1,
  iter = 50000,
  warmup = 10000,
  thin= 10,
  seed = 1919,
  init = inits,
  control = list(max_treedepth = 20, adapt_delta = 0.99))

#-----#
# Diagnostic
coda_chain_dummy <- As.mcmc.list(model_dummy, pars = c("kappa", "mu"))
summary(coda_chain_dummy)
acfplot(coda_chain_dummy)
geweke.diag(coda_chain_dummy, frac1=0.1, frac2=0.5)

#-----#
# What about the chains
rstan::traceplot(model_drugs, pars = "beta", inc_warmup = TRUE)
rstan::traceplot(model_drugs, pars = "mu", inc_warmup = TRUE)
rstan::traceplot(model_drugs, pars = "lambda", inc_warmup = TRUE)
rstan::traceplot(model_drugs, pars = "kappa", inc_warmup = TRUE)

#-----#
# Final iterations: 4000
# Posterior mean:
get_posterior_mean(model_drugs, pars = c("beta", "mu", "lambda", "kappa"))

#-----#
# What about the chains
rstan::traceplot(model_drugs, pars = "beta", inc_warmup = TRUE)
rstan::traceplot(model_drugs, pars = "mu", inc_warmup = TRUE)
rstan::traceplot(model_drugs, pars = "lambda", inc_warmup = TRUE)
rstan::traceplot(model_drugs, pars = "kappa", inc_warmup = TRUE)

#-----#
# Diagnostic
coda_chain_drugs <- As.mcmc.list(model_drugs, pars = c("kappa", "mu"))
summary(coda_chain_drugs)
acfplot(coda_chain_drugs)
geweke.diag(coda_chain_drugs, frac1=0.1, frac2=0.5)

#-----#

```

27/11/2020

919

```

# Random effects
plot(model_drugs, ask = T, pars = c("mu"))
plot(model_drugs, ask = T, pars = c("beta"))

# Diagnostic
coda_chain <- As.mcmc.list(model_drugs, pars = c("beta", "mu", "lambda", "kappa"))

acfplot(coda_chain)
geweke.diag(coda_chain, frac=0.1, frac2=0.5)

# Posterior
plot_post <- model_drugs %>%
  rstan::extract(c("beta"))
as.data.frame() %>%
  map_df(as_data_frame, .id = 'param')
cv <- rep(c("Age", "Length", "Age", "Length", "Age", "Length"), each = 40000)
gr <- rep(c("gr 1", "gr 2", "gr 3", "gr 4"), each = 8000)
data.frame(plot_post, cv = cv, gr = gr) %%

data.frame(plot_post, cv = cv, gr = gr) %%

ggplot(aes(value, fill = param)) +
  geom_density() +
  facet_grid(gr~cv, scales = 'free') +
  theme_minimal() +
  theme(legend.position="none")

plot_post <- model_drugs %>%
  rstan::extract(c("mu")) %>%
  as.data.frame() %>%
  map_df(as_data_frame, .id = 'param')

plot_post %>%
  ggplot(aes(value, fill = param)) +
  facet_wrap(~param, scales = 'free', ncol = 4) +
  theme_minimal() +
  xlim(c(4,6.5)) +
  theme(legend.position="none")

# model comparison

dati_list = list(N = nrow(data_drugs),
                 P = 2,
                 time = round(exp(data_drugs$time)),
                 cens = 1 - data.drugs$cens,
                 X = data_drugs[,4:5])

# STAN model

model_drugs_fix = stan(file = "Drugs_fix.stan",
                        data = dati_list,
                        chains = 1,
                        iter = 50000,
                        warmup = 10000,
                        thin= 10,
                        seed = 1919,
                        control = list(max_treedepth = 20, adapt_delta = 0.99))

# save(model_drugs_fix, file="model_drugs_fix.dat")
load("model_drugs_fix.dat")

# STAN model

model_drugs_fix = stan(file = "Drugs_fix.stan",
                        data = dati_list,
                        chains = 1,
                        iter = 50000,
                        warmup = 10000,
                        thin= 10,
                        seed = 1919,
                        control = 11, list(max_treedepth = 20, adapt_delta = 0.99))

# Diagnostic
coda_chain <- As.mcmc.list(model_drugs_fix, pars = c("beta", "mu", "lambda", "kappa"))

acfplot(coda_chain)
geweke.diag(coda_chain, frac=0.1, frac2=0.5)

# STAN model

model_drugs_dummy = stan(file = "Drugs_dummy.stan",
                           data = dati_list,
                           chains = 1,
                           iter = 50000,
                           warmup = 10000,
                           thin= 10,
                           seed = 1919,
                           control = list(max_treedepth = 20, adapt_delta = 0.99))

# save(model_drugs_dummy, file="model_drugs_dummy.dat")
load("model_drugs_dummy.dat")

# Diagnostic
rstan::traceplot(model_drugs_dummy, pars = c("mu", "inc_warmup = TRUE"))
rstan::traceplot(model_drugs_dummy, pars = c("lambda", "inc_warmup = TRUE"))
rstan::traceplot(model_drugs_dummy, pars = c("kappa", "inc_warmup = TRUE"))

# MODEL COMPARISON - WAIC
my_WAIC(model_drugs, "log_lik")
my_WAIC(model_drugs_fix, "log_lik")
summary(coda_chain)
acfplot(coda_chain)
geweke.diag(coda_chain, frac=0.1, frac2=0.5)

# Diagnostic
coda_chain <- As.mcmc.list(model_drugs_dummy, pars = c("mu", "lambda", "kappa"))

acfplot(coda_chain)
geweke.diag(coda_chain, frac=0.1, frac2=0.5)

```

BNP1

Riccardo Corradin
December 4, 2020

1/30

de Finetti theorem and Pólya Urn

de Finetti Th. & Pólya urn

Let \mathbb{X} be a Polish space, with \mathcal{X} denoting its Borel σ -field. Let $\{X_1, \dots, X_n\}$ be a finite sequence of \mathbb{X} -valued random elements defined on a common probability space $(\Omega, \mathcal{A}, \mathbb{P})$. We further assume the sequence exchangeable, in the sense that the distribution law of $\{X_1, \dots, X_n\}$ is invariant w.r.t. any permutation of the elements $\sigma : \mathbb{N}_n \rightarrow \mathbb{N}_n$.

$$\begin{array}{ccc} \mathbb{N}_4 & \xrightarrow{\sigma} & \mathbb{N}_4 \\ A_1 = \begin{bmatrix} \{X_1\} \\ \{X_2\} \end{bmatrix} & \xrightarrow{\quad} & \begin{bmatrix} \{X_1\} \\ \{X_4\} \end{bmatrix} = \sigma(A_1) \\ A_2 = [\{X_3\}] & \searrow & [\{X_2\}] = \sigma(A_2) \\ A_3 = [\{X_4\}] & \swarrow & [\{X_3\}] = \sigma(A_3) \end{array}$$

$$P(\{X_1, X_2\}, \{X_3\}, \{X_4\}) = P(\{X_1, X_4\}, \{X_2\}, \{X_3\})$$

2/30

de Finetti Th. & Pólya urn

We have a characterization of an exchangeable sequence thanks to the de Finetti representation theorem.

Theorem (de Finetti, 1937)

The sequence $\{X_1, \dots, X_n\}$ is exchangeable if and only if there exists a probability measure Q such that

$$P[X_1 = x_1, \dots, X_n = x_n] = \int \prod_{i=1}^n P[X_i = x_i | \theta] Q(d\theta)$$

The theorem states that a sequence is exchangeable iff it is conditionally independent.

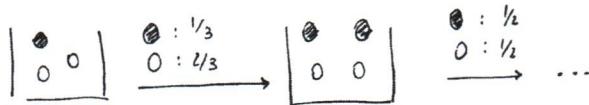
The measure Q is termed *de Finetti measure*, and it plays the role of prior distribution in Bayesian statistics.

is playing the role of prior distribution in Bayesian statistic

3/30

de Finetti Th. & Pólya urn

We consider a Pólya urn sampling scheme. Suppose we have an urn with balls of two 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 replace the ball in the urn with a new one of the same color.



We denote by:

- B the initial number of black balls;
- W the initial number of white balls;
- X_i the i -th Bernoulli sampling, where $X_i = \begin{cases} 1 : \text{black ball} \\ 0 : \text{white ball} \end{cases}$;
- S_n the number of black balls in n trials;
- $n - S_n$ the number of white balls in n trials;

4/30

de Finetti Th. & Pólya urn

Notice that a sequence sampled from a Pólya urn is exchangeable. We have:

$$P[X_4 = \{1, 0, 1, 1\}] = \frac{B}{B+W} \frac{W}{B+W+1} \frac{B+1}{B+W+2} \frac{B+2}{B+W+3}$$

$$P[X_4 = \{1, 1, 0, 1\}] = \frac{B}{B+W} \frac{B+1}{B+W+1} \frac{W}{B+W+2} \frac{B+2}{B+W+3}.$$

and the two coincides. More in general, for a n -sized sample from a Pólya urn scheme with two distinct color we have

$$P[X_1, \dots, X_n] = \frac{\Gamma(B + S_n)\Gamma(W + n - S_n)\Gamma(B + W)}{\Gamma(B)\Gamma(W)\Gamma(B + W + n)}$$

5/30

de Finetti Th. & Pólya urn

From the de Finetti theorem we have

$$P[X_1, \dots, X_n] = \int \prod_{i=1}^n \theta^{X_i} (1-\theta)^{1-X_i} Q(d\theta) = \int \theta^{S_n} (1-\theta)^{n-S_n} Q(d\theta)$$

Further, the measure Q is the cumulative distribution function of the limit frequency, which for the model we are considering corresponds to the cdf of

$$Y = \lim_{n \rightarrow \infty} \frac{1}{n} \sum_{i=1}^n X_i \sim \text{Beta}(B, W)$$

We can prove now that for an exchangeable sequence sampled from a Pólya urn sampling scheme, the de Finetti representation theorem holds.

6/30

de Finetti Th. & Pólya urn

proof.

We have

$$\begin{aligned} P[X_1, \dots, X_n] &= \int \theta^{S_n} (1-\theta)^{n-S_n} Q(d\theta) \\ &= \int \theta^{S_n} (1-\theta)^{n-S_n} \frac{\theta^{B-1} (1-\theta)^{W-1}}{B(B, W)} d\theta \\ &= \frac{B(B + S_n, W + n - S_n)}{B(B, W)} \int \frac{\theta^{B+S_n-1} (1-\theta)^{W+n-S_n-1}}{B(B + S_n, W + n - S_n)} d\theta \\ &= \frac{\Gamma(B + S_n)\Gamma(W + n - S_n)\Gamma(B + W)}{\Gamma(B)\Gamma(W)\Gamma(B + W + n)}, \end{aligned}$$

we substitute to Q the corresponding measure that we stated for the limit frequency (in our case we're measuring w.r.t. a Beta) ← probability associated to an exchangeable sequence produced from a polyà urn

which corresponds to the quantity of interest.

For an exchangeable sequence sampled from a Pólya urn the de Finetti theorem holds, and the de Finetti measure Q playing the role of prior distribution corresponds to a Beta distribution.

7/30

de Finetti Th. & Pólya urn

We can generalize the de Finetti representation theorem for a scenario with a infinite sequence and an infinite number of features. Let $X^{(\infty)}$ be an infinite sequence of observations, where the i -th element X_i of the sequence is taking values on a measurable space $(\mathbb{X}, \mathcal{X})$.

The sequence $X^{(\infty)}$ is exchangeable, for any $n \geq 1$ and any permutation σ of $\{1, \dots, n\}$, (X_1, \dots, X_n) coincides in distribution with $(X_{\sigma(1)}, \dots, X_{\sigma(n)})$.

Let $\mathbb{X}^{(\infty)} = \mathbb{X} \times \mathbb{X} \times \dots$ and \mathbb{P} be the space of probability measures on $(\mathbb{X}, \mathcal{X})$.

Theorem (de Finetti, 1937)

The sequence $X^{(\infty)}$ is exchangeable if and only if there exists a probability measure Q on \mathbb{P} such that, for any $n \geq 1$, and $A = A_1 \times A_2 \times \dots \times A_n \times \mathbb{X}^{(\infty)}$, with $A_i \in \mathcal{X}$, one has

$$\mathbb{P}[X^{(\infty)} \in A] = \int_{\mathbb{P}} \prod_{i=1}^n p(A_i) Q(dp).$$

8/30

An introduction to the DP

DP - Introduction

distribution over distributions : distribution taking values on a space of probability distributions

Let \mathbb{X} be a Polish space, with \mathcal{X} denoting its Borel σ -field. We further denote by \mathbb{P} the space of probability measures with support \mathbb{X} , endowed with its Borel σ -field \mathcal{P} .

Definition (Dirichlet process – Ferguson, 1973)

Let $\alpha(\cdot)$ be a finite, non-negative and diffuse measure defined on $(\mathbb{X}, \mathcal{X})$, with $\alpha(\mathbb{X}) < +\infty$. We say that P is distributed as a Dirichlet process with parameter $\alpha(\cdot)$ if for any partition $\{A_1, \dots, A_k\}$ of \mathbb{X} we have

$$(P(A_1), \dots, P(A_k)) \sim Dir(\alpha(A_1), \dots, \alpha(A_k)).$$

Notice that P is a random probability measure with support \mathbb{P} . We can further decompose the parameter of the DP in two terms $\alpha(\cdot) = aP_0(\cdot)$, with $P_0(\cdot) = \alpha(\cdot)/a$, which is termed base measure, and $a = \alpha(\mathbb{X})$, the strength / total mass parameter. Two different notations:

- $P \sim DP(\alpha(\cdot))$
- $P \sim DP(a, P_0(\cdot))$

9/30

DP - Moments

Let A be a measurable subset of \mathbb{X} , and $A^c = \mathbb{X} \setminus A$. Thanks to the definition of DP we have that

$$P(A) \sim Beta(\alpha(A), \alpha(A^c))$$

then we have:

- $E[P(A)] = \frac{\alpha(A)}{\alpha(A)+\alpha(A^c)} = P_0(A)$
- $Var(P(A)) = \frac{\alpha(A)\alpha(A^c)}{(\alpha(A)+\alpha(A^c))^2(\alpha(A)+\alpha(A^c)+1)} = \frac{a^2 P_0(A) P_0(A^c)}{a^2 (P_0(A) + P_0(A^c))(a+1)} = \frac{P_0(A) P_0(A^c)}{(a+1)}$

a plays the role of a precision parameter.

The Dirichlet process is a.s. discrete random probability measure with

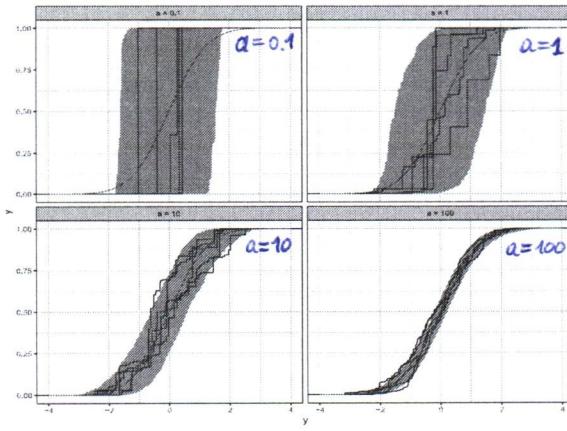
$$P(dt) = \sum_{j=1}^{\infty} W_j \delta_{X_j^*}(dt)$$

– double randomness inside the process: both $\{X_j^*\}_{j \geq 1}$ and $\{W_j\}_{j \geq 1}$ are random

with $\delta_{X_j^*}(\cdot)$ a Dirac measure in X_j^* , $\{W_j\}_{j \geq 1}$ a sequence of random weights, $\{X_j^*\}_{j \geq 1}$ a sequence of random elements i.i.d. from $P_0(\cdot)$, with $\{W_j\}_{j \geq 1} \perp \!\!\! \perp \{X_j^*\}_{j \geq 1}$

10/30

DP - Moments



11/30

DP - Posterior

Suppose we collected an exchangeable sample X_1, \dots, X_n from a DP, showing X_1^*, \dots, X_k^* distinct values with frequencies n_1, \dots, n_k , i.e.

$$X_1, \dots, X_n \mid P \stackrel{\text{iid}}{\sim} P, \quad P \sim DP(\alpha(\cdot)),$$

and we are interested into characterize the posterior distribution of P , given X_1, \dots, X_n .

The posterior distribution is again a Dirichlet process, and is given by

$$P \mid X_1, \dots, X_n \sim DP \left(\alpha(\cdot) + \sum_{j=1}^k n_j \delta_{X_j^*}(\cdot) \right)$$

or equivalently $P \mid X_1, \dots, X_n \sim DP \left(aP_0(\cdot) + \sum_{j=1}^k n_j \delta_{X_j^*}(\cdot) \right)$

► posterior total mass

$$\int_X aP_0(x) + \sum_{j=1}^k n_j \delta_{X_j^*}(x) dx = a + n$$

► posterior expectation on a measurable set $A \subseteq \mathbb{X}$

$$\mathbb{E}[P(A) \mid X_1, \dots, X_n] = \frac{a}{a+n} P_0(A) + \frac{n}{a+n} \sum_{j=1}^k \frac{n_j}{n} \delta_{X_j^*}(A)$$

12/30

DP - Number of unique elements (aka number of clusters)

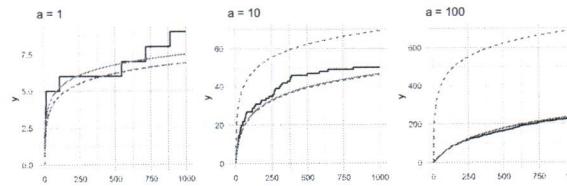
Let K_n be the number of unique elements out of a sample of size n . The expected number of elements out of a sample of size n is equal to

$$\mathbb{E}[K_n] = \sum_{j=1}^n \frac{a}{a+j-1} = a(\psi(a+n) - \psi(a)) \simeq a \log \left(\frac{a+n}{a} \right) \simeq a \log(n)$$

and

$$\text{Var}(K_n) = a(\psi(a+n) - \psi(a)) + a(\psi'(a+n) - \psi'(a)) \simeq a \log \left(\frac{a+n}{a} \right)$$

for $n \rightarrow \infty$, with $n > a \gg 0$ and $\psi(\cdot)$ denoting the digamma function.



13/30

DP - Self similarity

What happen when we restrict the DP to a subset of its support? For example, what is the distribution of $P|_B$ with $B \subset \mathbb{X}$?

Theorem (Self similarity)

Let $P \sim DP(a, P_0(\cdot))$ a Dirichlet process with support \mathbb{X} , with $P_0(\cdot)$ a diffuse probability measure on \mathbb{X} . Let B a measurable subset of \mathbb{X} , and $B^c = \mathbb{X} \setminus B$. Then

$$P(B), \quad P(B^c), \quad P|_B, \quad P|_{B^c}$$

are mutually independent, for any $B : P_0(B) > 0$. Further

$$P|_B \sim DP(aP_0(B), P_0(\cdot)|_B)$$

$$P|_{B^c} \sim DP(aP_0(B^c), P_0(\cdot)|_{B^c})$$



A Dirichlet process restricted on a subset of its support is still distributed as a Dirichlet process, and is independent on what happens on the rest of the support.

14/30

DP - Self similarity

[PROOF] Consider a partition A_1, \dots, A_m of B , and a partition C_1, \dots, C_q of B^c . Then we have that over the two partitions joint together

$$\begin{aligned} \mathbf{X} &= (P(A_1), \dots, P(A_m), P(C_1), \dots, P(C_q)) \\ &\sim Dir(\alpha(A_1), \dots, \alpha(A_m), \alpha(C_1), \dots, \alpha(C_q)) \end{aligned}$$

Moreover, let

$$Z_1 = \sum_{j=1}^m X_j, \quad Z_2 = \sum_{j=m+1}^{m+q} X_j, \quad \mathbf{Y}_1 = \left(\frac{X_1}{Z_1}, \dots, \frac{X_m}{Z_1} \right), \quad \mathbf{Y}_2 = \left(\frac{X_{m+1}}{Z_2}, \dots, \frac{X_{m+q}}{Z_2} \right)$$

Thanks to the aggregation property of the Dirichlet distribution, the previous quantities are mutually independent. Further, \mathbf{Y}_1 and \mathbf{Y}_2 are defining two distinct Dirichlet process on B and B^c , with base measure $P_0(\cdot) |_B$ and $P_0(\cdot) |_{B^c}$ respectively, and total mass $aP_0(B)$ and $aP_0(B^c)$.

15/30

Chinese restaurant process

= process of allocation of elements to different groups

Suppose we have an infinite queue in front of a restaurant, with an infinite number of tables.

We are interested into modeling the table allocation of the customers, as far as they continue to enter in the restaurant.

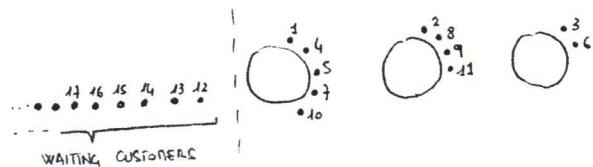
Conditionally on the allocation of the first n customers, the next $n+1$ customer can sit to an already occupied table with probability proportional to the current number of customers at that table, or at a new table.

- For the j -th table denoting by n_j the current number of customers at table j , the probability that the $n+1$ customer will sit at the table j is proportional to n_j .
- The probability that the $n+1$ customer will sit at a new table is proportional to a .

16/30

Chinese restaurant process

"La Riccardo" RESTAURANT



The $n+1 = 12$ customer can sit at the first table with probability proportional to $n_1 = 5$, at the second table with probability proportional to $n_2 = 4$, at the third table with probability proportional to $n_3 = 2$, or at a new table with probability proportional to a .

i.e. we have $p_1 = \frac{n_1}{a+n_1+n_2+n_3} = \frac{5}{a+11}$,

and $p_2 = \frac{4}{a+11}$, $p_3 = \frac{2}{a+11}$, and $p_{new} = \frac{a}{a+11}$.

17/30

Chinese restaurant process

Let Z_i denotes the allocation of the i -th customer. We then have

$$P[Z_{n+1} = t | Z_1, \dots, Z_n] = \frac{a}{a+n} \delta_{k+1}(t) + \frac{n}{a+n} \sum_{j=1}^k \frac{n_j}{n} \delta_j(t)$$

Where k is the number of filled table by the first n customers; and n_j is the number of customers at the j -th table.

Remarks:

- CRP shows the "rich-gets-richer" phenomenon, where large clusters grow larger faster (the number of clusters grows logarithmic w.r.t. the number of observations).
- We are not interested in any specific feature characterizing each group;
- The allocation probability is in relation with a Dirichlet process;
- The probability associated to a specific allocation is

$$P(Z_1, \dots, Z_n) = \frac{a^k}{(a)_n} \prod_{j=1}^k (n_j - 1)!$$

with $(a)_b$ denoting the rising factorial.

probability of allocating the new customer to a new table

probability of allocating the new customer to one of the occupied tables (NOTICE: we're not making distinctions between groups)

18/30

Generalized Pólya urn

Suppose that we are now interested also in some features characterizing each group.

For example, we are in an Pólya urn regime. 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 .

But we have also a positive probability of sampling a ball of a new color, from another urn with an infinite number of balls $\alpha(\cdot)$, where each color is unique.

probability of
sampling a specific
color
(or generally a
specific value)

Given a sample X_1, \dots, X_n , with X_1^*, \dots, X_k^* unique elements and n_1, \dots, n_k frequencies, we have

$$P[X_{n+1} \in dt | X_1, \dots, X_n] \propto \alpha(dt) + \sum_{j=1}^k n_j \delta_{X_j^*}(dt)$$

The previous reminds the predictive distribution of a Dirichlet process.
Early studied by Blackwell and MacQueen.

in the case of Chinese restaurant we were not interested in group's (table's) features and the characteristics

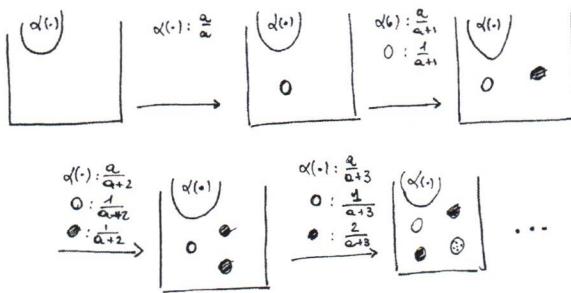
$P(\text{new color})$

$P(\text{one of the already-there colors})$

Generalization of the Pólya urn

19/30

Generalized Pólya urn



How the generalized Pólya urn works in practice. At each step we can sample an already observed color, or a new one from $\alpha(\cdot)$.

20/30

Generalized Pólya urn

Remarks:

► Then the joint law of (X_1, \dots, X_n) can be decomposed as

$$\begin{aligned} P(X_1, \dots, X_n) &= P(X_n | X_{n-1}, \dots, X_1) P(X_{n-1} | X_{n-2}, \dots, X_1) \dots P(X_1) \\ &= \prod_{i=1}^n \left[\frac{a}{a+i-1} P_0(\cdot) + \frac{i-1}{a+i-1} \sum_{j=1}^{k_{[i]}} \frac{n_{j[i]}}{i-1} \delta_{X_{j[i]}^*} \right] \end{aligned}$$

with $k_{[i]}$ denoting the number of already observed distinct elements at the i -th sampling, $X_{1[i]}^*, \dots, X_{k_{[i]}[i]}^*$ the unique values, and $n_{1[i]}, \dots, n_{k_{[i]}[i]}$ their frequencies;

► The sampled sequence is exchangeable, i.e.

$$P(X_1, \dots, X_n) = P(X_{\sigma(1)}, \dots, X_{\sigma(n)})$$

with $\sigma : \mathbb{N}_n \rightarrow \mathbb{N}_n$ a permutation of $1, \dots, n$;

► By sampling from a generalized Pólya urn, we obtain a trajectory from a Dirichlet process;

► It can be viewed as a CRP where we assign to each table an unique dish chosen from P_0 .

21/30

Stick-breaking

We want to represent $P = \sum_{j=1}^{\infty} W_j \delta_{X_j^*}$ in a useful way. First, we recall that $\{W_j\}_{j \geq 1}$ are independent of $\{X_j^*\}_{j \geq 1}$.

We know that $\{X_j^*\}_{j \geq 1}$ are i.i.d from $P_0(\cdot)$.

We need to characterize the distribution of $\{W_j\}_{j \geq 1}$. Possible solution: stick-breaking construction: instead of characterizing the value of W_j , we focus on the distribution of the mass that we are stealing at the j -th step from the residual mass after $j-1$ steps.

Suppose we are stealing mass β at the first step on the atom X_1^* . From the posterior distribution of a DP we have

$$P | X_1^* \sim DP \left(a+1, \frac{aP_0(\cdot) + \delta_{X_1^*}(\cdot)}{a+1} \right)$$

We further know that defining a partition $(X_1^*, \mathbb{X} \setminus X_1^*)$ of \mathbb{X} we have

$$\begin{aligned} (P(X_1^*), P(\mathbb{X} \setminus X_1^*)) &\sim Dir \left((a+1) \frac{aP_0(X_1^*) + \delta_{X_1^*}(X_1^*)}{a+1}, (a+1) \frac{aP_0(\mathbb{X}) + \delta_{X_1^*}(\mathbb{X})}{a+1} \right) \\ &\sim Dir(1, a) \end{aligned}$$

22/30

Stick-breaking

We then have that P has a point mass located in X_1^* , and can be rewritten as

$$G = V_1 \delta_{X_1^*} + (1 - V_1)P|_{X \setminus X_1^*}$$

Thanks to the self similarity, we know that $P|_{X \setminus X_1^*}$ is still distributed as a DP with $P_0(\cdot)|_{X \setminus X_1^*} \stackrel{\text{def}}{=} P_0(\cdot)$, with mass a . Further, we know the marginal distribution for a dimension of a Dirichlet random variable, and we have that

$$V_1 \sim \text{Beta}(1, a)$$

We then have

$$P \sim DP(a, P_0(\cdot))$$

$$P \sim V_1 \delta_{X_1^*} + (1 - V_1)P_1, \quad \text{with } P_1 \sim DP(a, P_0(\cdot))$$

$$P \sim V_2 \delta_{X_2^*} + (1 - V_1)(V_2 \delta_{X_2^*} + (1 - V_2)P_2), \quad \text{with } P_2 \sim DP(a, P_0(\cdot))$$

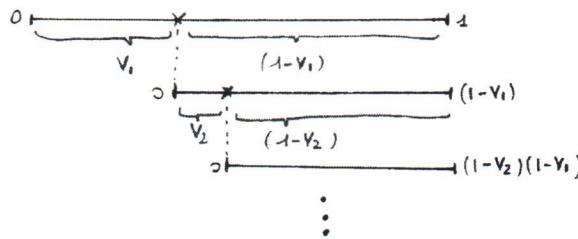
⋮

$$P = \sum_{j=1}^{\infty} W_j \delta_{X_j^*}$$

with $\{X_j^*\}_{j \geq 1}$ i.i.d from $P_0(\cdot)$, $W_j = V_j \prod_{\ell < j} (1 - V_\ell)$ and $V_j \sim \text{Beta}(1, a)$.

23/30

Stick-breaking



How the stick-breaking works in practice with the distribution of the weights. Each V_j is distributed according to a $\text{Beta}(1, a)$ distribution.

24/30

Stick-breaking

The stick-breaking construction is more general than the Dirichlet process case.

We need some regularity conditions:

$$\sum_{j=1}^{\infty} W_j = 1 \text{ a.s.}, \quad \text{IFF} \quad \sum_{j=1}^{\infty} \mathbb{E}[\log(1 - V_j)] = -\infty$$

We are dealing with a sequence $\{V_j\}_{j \geq 1} \stackrel{\text{iid}}{\sim} \text{Beta}(1, a)$.

It is sufficient that $P[V_1 > 0] > 0$.

The posterior distribution can also be represented in terms of stick-breaking construction, with updated weights.

$$\begin{aligned} P|_{X_1, \dots, X_n} &= \sum_{j=1}^{\infty} W_j \delta_{X_j^*}(\cdot) \\ W_j &= V_j \prod_{\ell < j} (1 - V_\ell) \\ V_j &\stackrel{\text{ind}}{\sim} \text{Beta}(1 + n_j, a + m_j) \end{aligned}$$

where $n_j = \sum_{i=1}^n \mathbb{I}_{[X_i = X_j^*]}$ and $m_j = n - \sum_{\ell=1}^j \sum_{i=1}^n \mathbb{I}_{[X_i = X_\ell^*]}$

25/30

Sampling from a Dirichlet process

we can sample a trajectory or a realization from the random prob. measure

Sampling from a DP: a trajectory

To sample a trajectory of size n from a Dirichlet process, we can exploit the Pólya urn sampling scheme. We start from X_1 where

$$P[X_1 \in dt] = \frac{a}{a} P_0(dt) \sim P_0(dt)$$

then we have

$$\begin{aligned} P[X_2 \in dt | X_1] &= \frac{a}{a+1} P_0(dt) + \frac{1}{a+1} \delta_{X_1}(dt) \\ P[X_3 \in dt | X_1, X_2] &= \frac{a}{a+2} P_0(dt) + \frac{2}{a+2} \sum_{j=1}^k \frac{n_j}{2} \delta_{X_j^*} \\ P[X_n \in dt | X_1, \dots, X_{n-1}] &= \frac{a}{a+n-1} P_0(dt) + \frac{n-1}{a+n-1} \sum_{j=1}^k \frac{n_j}{n-1} \delta_{X_j^*} \end{aligned}$$

and the sequence produced X_1, \dots, X_n is an n -size sample from a Dirichlet process with mass parameter a and base measure P_0 .

26/30

Sampling from a DP: updating a trajectory

Let $\mathbf{X}_{(i)} = \{X_1, \dots, X_n\} \setminus X_i$. Let $\mathcal{L}(Y)$ denotes the distributional law of Y . Then we have

$$\begin{aligned} \mathcal{L}(X_i | \mathbf{X}_{(i)}) &= \int_{P_X} \mathcal{L}(X_i, P | \mathbf{X}_{(i)}) dP \\ &= \int_{P_X} \underbrace{\mathcal{L}(X_i | P, \mathbf{X}_{(i)})}_{X_i | P \sim P} \underbrace{\mathcal{L}(P | \mathbf{X}_{(i)})}_{P | \mathbf{X}_{(i)} \sim DP(a, 1 + \sum_{j \neq i} n_j \delta_{X_j})} dP \\ &= \int_{P_X} P \mathcal{L}(dP | \mathbf{X}_{(i)}) = \mathbb{E}[P | \mathbf{X}_{(i)}] \\ &= \frac{a}{a+n-1} P_0(\cdot) + \frac{n-1}{a+n-1} \sum_{j=1}^{k_{(i)}} \frac{n_{j(i)}}{n-1} \delta_{X_{j(i)}^*} \quad (\text{OK}) \end{aligned}$$

where $k_{(i)}$ are the number of distinct values in $\mathbf{X}_{(i)}$, $X_{j(i)}^*, \dots, X_{k_{(i)}(i)}^*$ denote the distinct values out of $\mathbf{X}_{(i)}$, and $n_{j(i)}$ the frequency of the j -th distinct value.

we can update a trajectory updating element by element conditional to the rest of the sample and sampling for each element something that we already observed (OK) or something new (OK) with probabilities (OK) and (OK)

27/30

Sampling from a DP: a truncated realization

We might be interested into sampling a realization from a Dirichlet process.

WARNING: we have an infinite number of dimensions.

A first attempt: we can truncate the stick-breaking representation.

$$P(dt) \simeq \tilde{P}(dt) = \sum_{j=1}^J W_j \delta_{X_j^*}(dt)$$

with

$$\begin{aligned} W_j &= V_j \prod_{\ell < j} (1 - V_\ell), \quad j = 1, \dots, J-1 \\ V_j &\sim \text{Beta}(1, a), \quad j = 1, \dots, J-1 \\ V_J &= 1 \end{aligned}$$

Notice that:

- we introduce an approximation;
- we can control in some way the approximation;
- we can use a random truncation, for example $J \sim \text{Poi}(\lambda)$.

28/30

Sampling from a DP: a sliced realization

Instead of truncating truncate the stick-breaking construction, we can augment the problem and implement a slice sampling strategy, moving from $\mathcal{L}(P)$ to $\mathcal{L}(P, U)$, with $U \sim \text{Unif}(0, 1)$, such that

$$\mathcal{L}(P, U) = \sum_{j=1}^{\infty} \mathbb{I}_{[U \leq W_j]} \delta_{X_j^*}$$

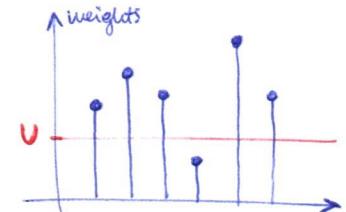
Notice that conditionally on a specific value for U , we can represent P with a finite number of component:

- Only the atoms with $W_j \geq U$ are included.
- Exists M such that $\sum_{j=M}^{\infty} W_j < U$, so that there are no weights after M which can be possibly larger than U .

Further

$$\int \mathcal{L}(P, U) du = \sum_{j=1}^{\infty} \int \mathbb{I}_{[U \leq W_j]} \delta_{X_j^*} \mathbb{I}_{[0,1]}(u) du = \sum_{j=1}^{\infty} W_j \delta_{X_j^*}$$

If we marginalize wrt U , we obtain $\mathcal{L}(P)$.



Conditionally on the augmented variable $U \sim \text{Unif}(0, 1)$, we select just a finite number of weights, all the W_j 's s.t. $W_j > U$

29/30

4/12/2020
5/6

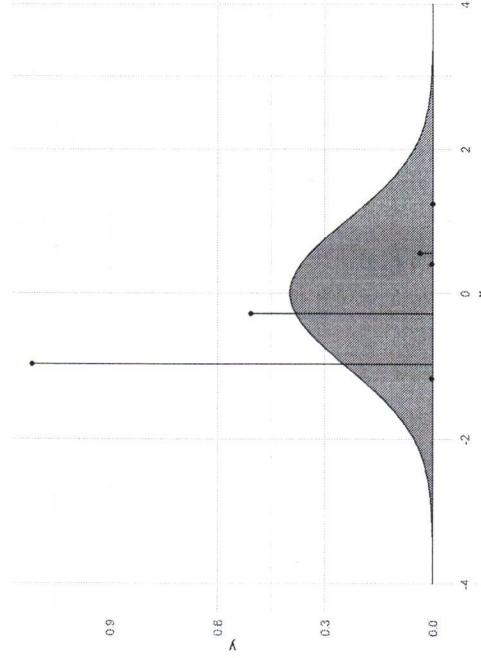
```
PA_true * (1 - PA_true) / (a + 1)
```

```
## [1] 0.07424874
```

```
sample[[1]]
```

```
## frequency value
## [1,] 3841 -0.276878
## [2,] 6691 -0.9641521
## [3,] 224 0.5612653
## [4,] 23 0.4112657
## [5,] 17 -1.1689911
## [6,] 4 1.2365513
```

```
data.frame(x = sample[[1]][1], y = sample[[1]][1]/60000) %>%
  ggplot() +
  geom_ribbon(data = data.frame(x = seq(-4, 4, by = 0.05),
                                 y = rep(0, 161),
                                 ymin = dnorm(seq(-4, 4, by = 0.05))),
               aes(x = x, ymin = ymin, ymax = ymax), fill = "orange") +
  geom_point(aes(x = x, y = y)) +
  theme_minimal() +
  geom_segment(aes(x = x, y = rep(0, length(x)), xend = x, yend = y))
```



```
#-- increase the mass parameter a = 25
a <- 25
sample2 <- traj_sample(100, 10000, a)
```

```
PA <- mean(unlist(
  lapply(sample2, function(x) sum(x[,2] < 1 && x[,2] > -2, 1)) / 10000
))

PA
```

```
## [1] 0.79
```

```
PA * (1 - PA) / (a + 1)
```

```
## [1] 0.00380769
```

```
## [1] 0.8155946
```

```
# R & STAN are friends!
library(RStan)
library(rstanarm)

# for plots
library(ggplot2)
library(tidyverse)
library(dplyr)

library(purrr)
library(BGcl)
require(BGcl)

require(BGpbor)

#-----#
# SAMPLING m TRAJECTORY FROM DP of Length n
# with Gaussian(θ, 1) base measure
# using the predictive distribution
#-----#
traj_sample <- function(m, n, alpha){
  result <- list()
  pb <- txtProgressBar(θ = m, style = 3)
  for(i in 1:m){
    temp_seq <- c()
    for(j in 1:n){
      index <- sample(1:(length(temp_seq) + 1), 1, prob = c(temp_seq, alpha))
      if(index <= length(temp_seq)){
        temp_seq[index] <- temp_seq[index] + 1
      } else {
        temp_seq[index] <- 1
      }
      temp_mat <- cbind(temp_seq, rnorm(length(temp_seq), θ, 1))
      colnames(temp_mat) <- c("frequency", "value")
      result[[j]] <- temp_mat
      setTxtProgressBar(pb, j)
    }
    close(pb)
    return(result)
  }
}

#-----#
# DIFFERENT TRIALS
# A = (-2, 1)
```

```
a <- 1
sample1 <- traj_sample(100, 10000, a)

PA <- mean(
  unlist(
    lapply(sample1, function(x) sum(x[,2] < 1 && x[,2] > -2, 1)) / 10000
  )
)

PA
```

```
## [1] 0.82
```

```
PA * (1 - PA) / (a + 1)
```

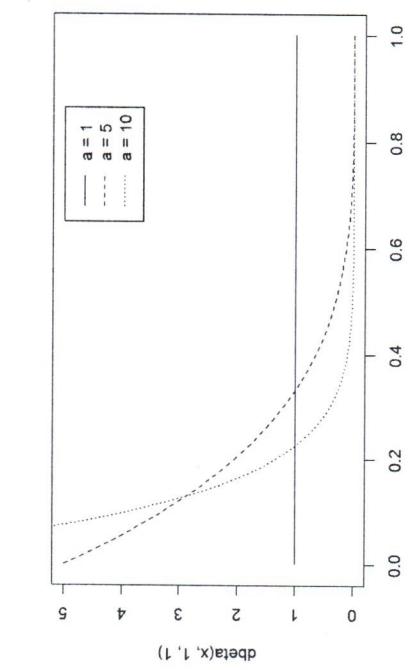
```
## [1] 0.0738
```

```
PA_true <- pnorm(1) - pnorm(-2)
PA_true
```

```
## [1] 0.8155946
```


4/12/2020

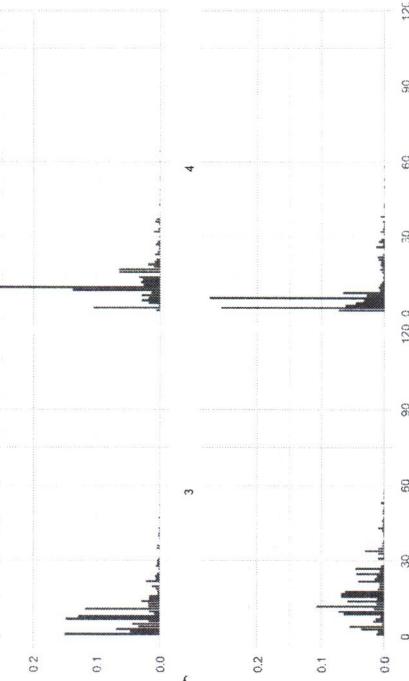
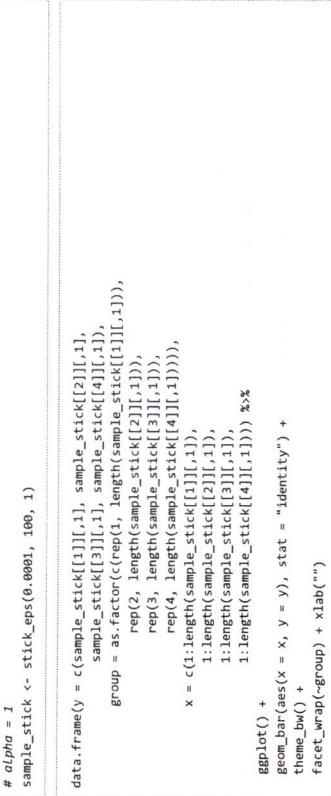
6/6



```
#-----#
# Sampling
stick_eps <- function(eps, m, alpha){
  result <- list()

  pb <- txtProgressBar(0, m, style = 3)
  for(j in 1:m){
    temp_seq_V <- rbeta(1,1,alpha)
    temp_seq_W <- temp_seq_V[1]
    k <- 2
    while(sum(temp_seq_W) < 1 - eps){
      temp_seq_V[k] <- rbeta(1,1,alpha)
      temp_seq_W[k] <- prod(1 - temp_seq_V[1:(k - 1)]) * temp_seq_V[k]
      k <- k + 1
    }
    result[[j]] <- cbind(temp_seq_W, rnorm(length(temp_seq_W), 0, 1))
    setTxtProgressBar(pb, j)
  }
  close(pb)
  return(result)
}

# alpha = 1
sample_stick <- stick_eps(0.0001, 100, 1)
```



```
# alpha = 10
sample_stick <- stick_eps(0.0001, 100, 10)

data.frame(y = c(sample_stick[[1]][1], sample_stick[[2]][1]),
           sample_stick[[3]][1], sample_stick[[4]][1]),
           group = as.factor(rep(1, length(sample_stick[[1]][1])),
                           rep(2, length(sample_stick[[2]][1])),
                           rep(3, length(sample_stick[[3]][1])),
                           rep(4, length(sample_stick[[4]][1])))),
           x = c(1:length(sample_stick[[1]][1]),
                1:length(sample_stick[[2]][1]),
                1:length(sample_stick[[3]][1]),
                1:length(sample_stick[[4]][1]))))

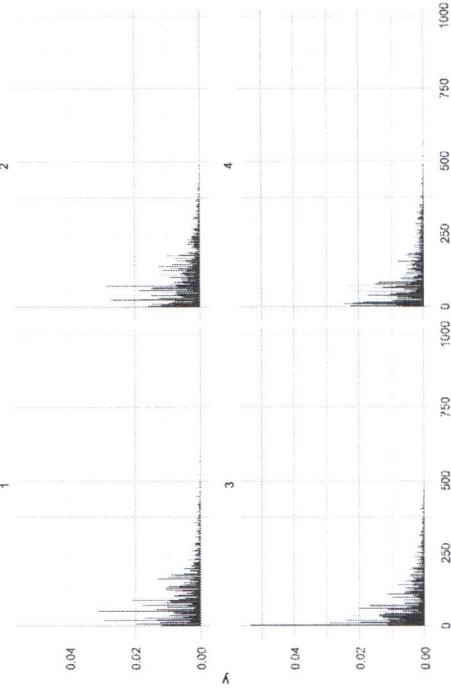
ggplot() +
  geom_bar(aes(x = x, y = y), stat = "identity") +
  theme_minimal() +
  facet_wrap(~group) + xlab("")
```

```
# alpha = 100
sample_stick <- stick_eps(0.0001, 100, 100)
```

```

data.frame(y = c(sample_stick[[1]][1], sample_stick[[2]][1], sample_stick[[3]][1], sample_stick[[4]][1]), group = as.factor(c(rep(1, length(sample_stick[[1]])), rep(2, length(sample_stick[[2]]))), rep(3, length(sample_stick[[3]]))), rep(4, length(sample_stick[[4]]))), x = c(1:length(sample_stick[[1]]), 1:length(sample_stick[[2]]), 1:length(sample_stick[[3]]), 1:length(sample_stick[[4]]))), p1 <- ggplot() + geom_bar(aes(x = x, y = y), stat = "identity") + theme_minimal() + facet_wrap(~group) + xlab("") + ggtitle("a = 1")
p2 <- ggplot() + geom_line(aes(x = x, y = y)) +
  theme_minimal() + xlab("") +
  geom_line(data = data.frame(x = 1:10^3, y = 1 * log(1:10^3)), mapping = aes(x = x, y = y), col = 2, lty = 2) +
  geom_line(data = data.frame(x = 1:10^3, y = 1 * log(1:10^3) / 100), mapping = aes(x = x, y = y), col = 4, lty = 2) +
  ggtitle("a = 10")
p3 <- ggplot() + geom_line(aes(x = x, y = y)) +
  theme_minimal() + xlab("") +
  geom_line(data = data.frame(x = 1:10^3, y = 10 * log(1:10^3)), mapping = aes(x = x, y = y), col = 2, lty = 2) +
  geom_line(data = data.frame(x = 1:10^3, y = 10 * log((1:10^3) / 100)), mapping = aes(x = x, y = y), col = 4, lty = 2) +
  ggtitle("a = 100")
ggarrange(p1, p2, p3, ncol = 3)

```



```

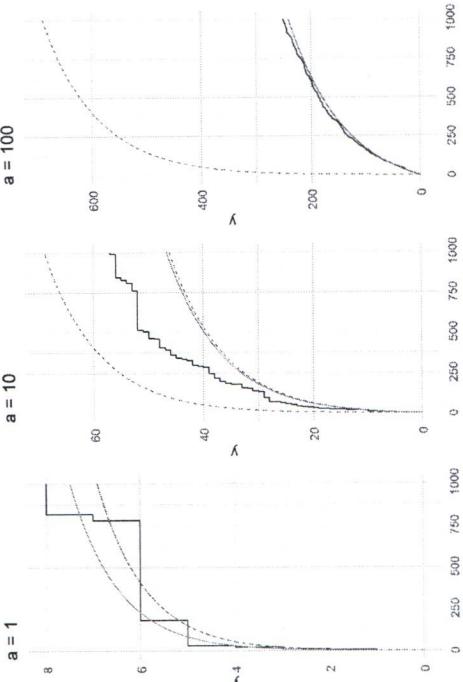
#-----#
# The number of clusters
traj_sample_unique <- function(n, alpha){
  result <- 1
  pb <- txtProgressBar(0, n, style = 3)
  for(i in 2:n){
    index <- sample(1:(length(unique(result))) + 1), 1, prob = c(table(result), alpha))
    result <- c(result, index)
    setTxtProgressBar(pb, i)
  }
  close(pb)
  return(result)
}

sample_one1 <- traj_sample_unique(10^3, 1)

sample_one2 <- traj_sample_unique(10^3, 10)

sample_one3 <- traj_sample_unique(10^3, 100)

```



BNP2

Riccardo Corradin
December 10, 2020

Recap : DIRICHLET PROCESS

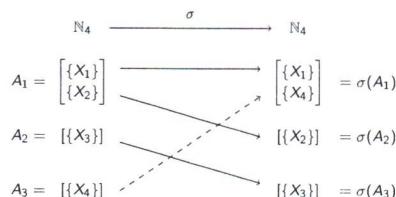
combinatorial process and its realizations are sequences of atoms which possibly show ties.
In general it's a probability measure with random jumps at random locations (we can have an infinite but countable jumps and locations)

1/26

Random partitions

Exchangeability

Let \mathbb{X} be a Polish space, with \mathcal{X} denoting its Borel σ -field. Let $\{X_1, \dots, X_n\}$ be a finite sequence of \mathbb{X} -valued random elements defined on a common probability space $(\Omega, \mathcal{A}, \mathbb{P})$. We further assume the sequence exchangeable, in the sense that the distribution law of $\{X_1, \dots, X_n\}$ is invariant w.r.t. any permutation of the elements $\sigma : \mathbb{N}_n \rightarrow \mathbb{N}_n$.



$$P(\{X_1, X_2\}, \{X_3\}, \{X_4\}) = P(\{X_1, X_4\}, \{X_2\}, \{X_3\})$$

2/26

Partitions

Let $\mathbb{N}_n = \{1, \dots, n\}$ space of first n integer.

Definition

We say that $\psi_n = \{A_1, \dots, A_k\}$ is a *partition* of \mathbb{N}_n if the sets $\{A_1, \dots, A_k\}$ are nonempty and such that $\mathbb{N}_n = \bigcup_i A_i$ and $A_i \cap A_j = \emptyset$ for all $i \neq j$.

Let \mathcal{B} be the space of all partitions of \mathbb{N}_n and \mathcal{B} its discrete σ -field.

Definition

Let $(\Omega, \mathcal{F}, \mathbb{P})$ be a generic probability space. A random partition is any measurable function $\Psi_n : (\Omega, \mathcal{F}, \mathbb{P}) \rightarrow (\mathcal{B}, \mathcal{B})$.

The definition of Ψ_n induces a probability measure on the space of partitions \mathcal{B} .

Given a partition $\psi_n = \{A_1, \dots, A_k\}$ of \mathbb{N}_n , the sequence of cardinalities of the blocks of ψ_n , that is $C_n = (n_1, \dots, n_k)$, with $n_i = |A_i|$, is named composition of n . A random partition Ψ_n induces a random composition C_n . Note that the space of all compositions of n is the simplex

$$\Delta_n^k = \left\{ (n_1, \dots, n_k) \in \mathbb{N}_n^k : \sum_{j=1}^k n_j = n, n_j > 0 \right\}.$$

3/26

Exchangeable partitions

Let σ be a permutation of \mathbb{N}_n and $\sigma(A) = \{\sigma(i) : i \in A\}$.

Definition (Exchangeable partition)

We say that a random partition Ψ_n is exchangeable if, given any permutation $\sigma : \mathbb{N}_n \rightarrow \mathbb{N}_n$, the distribution of Ψ_n is invariant with respect to σ , i.e. for any partition $\{A_1, \dots, A_k\}$ of \mathbb{N}_n and for any permutation σ , the following holds

$$P[\Psi_n = \{A_1, \dots, A_k\}] = P[\Psi_n = \{\sigma(A_1), \dots, \sigma(A_k)\}].$$

Equivalently, we say that a random partition Ψ_n is exchangeable if there exists a symmetric function $p_k^{(n)} : \mathcal{C}_n \rightarrow [0, 1]$ such that, for every partition $\{A_1, \dots, A_k\}$ of \mathbb{N}_n

$$P[\Psi_n = \{A_1, \dots, A_k\}] = p_k^{(n)}(|A_1|, \dots, |A_k|) = p_k^{(n)}(n_1, \dots, n_k).$$

The function $p_k^{(n)}$ is named exchangeable partition probability function (EPPF) of Ψ_n .

we can change the elements between the groups and also the labels, the only important thing is the cardinality of the groups (whatever the name of the group is, whatever the elements of the group are)

4/26

EPPFs

Remark: $p_k^{(n)}$ symmetric \Rightarrow the cluster labels are exchangeable.

probability of dividing one obs. in one group $\Rightarrow 1$

the probability is invariant w.r.t. permutations of the blocks

An EPPF satisfies the following properties:

- i) $p_1^{(1)} = 1$;
- ii) for any $(n_1, \dots, n_k) \in \Delta_n^k$, with $n \geq 1$ and $1 \leq k \leq n$,

$$p_k^{(n)}(n_1, \dots, n_k) = p_k^{(n)}(n_{\sigma(1)}, \dots, n_{\sigma(k)}),$$

where $\sigma(\cdot)$ is a permutation of $(1, \dots, k)$;

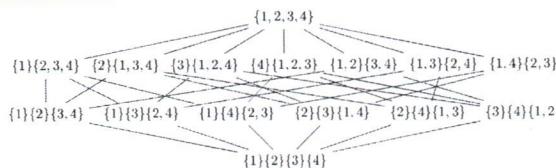
- iii) for any $(n_1, \dots, n_k) \in \Delta_n^k$, with $n \geq 1$ and $1 \leq k \leq n$, the following addition rule holds

$$p_k^{(n)}(n_1, \dots, n_k) = \sum_{j=1}^k p_k^{(n+1)}(n_1, \dots, n_j + 1, \dots, n_k) + p_{k+1}^{(n+1)}(n_1, \dots, n_k, 1). \quad (1)$$

5/26

EPPFs

Problem: the space of the possible partition of n elements is large, and is growing fast as far as n grows.



We have

$$S_{n,k} = \frac{1}{k!} \sum_{j=0}^k (-1)^j \binom{k}{j} (k-j)^n$$

(Stirling number of the second kind) ways to partition n elements in k groups, and

$$B_n = \sum_{k=1}^n S_{n,k}$$

(Bell number) ways to partition n elements. B_n is exploding as far as n grows.

6/26

Connection with the Dirichlet process

DP - EPPF

Let P be a random probability measure distributed as a Dirichlet process with total mass $a > 0$ and base measure $P_0(\cdot)$, and let X_1, \dots, X_n be an exchangeable sample from P , i.e.

$$X_1, \dots, X_n \mid P \stackrel{\text{iid}}{\sim} P$$

$$P \sim DP(a, P_0(\cdot))$$

Due to the a.s. discreteness of P , the sample can possibly show ties. We denote by X_1^*, \dots, X_k^* the unique values in X_1, \dots, X_n , with frequencies n_1, \dots, n_k , for $1 \leq k \leq n$.

TIES IN THE SEQUENCE \Leftrightarrow LATENT PARTITION

The expectation of the distribution of a sequence (X_1, \dots, X_n) sampled from a Dirichlet process

$$\begin{aligned} \mathbb{E}[\mathcal{L}(X_1, \dots, X_n)] &= \mathbb{E}\left[\prod_{i=1}^n P(X_i)\right] = \mathbb{E}\left[\prod_{j=1}^k P^{n_j}(X_j^*)\right] \\ &= \underbrace{p_k^{(n)}(n_1, \dots, n_k)}_{\text{distribution of the partition}} \underbrace{\prod_{j=1}^k P_0(X_j^*)}_{\text{distribution of the atoms}} \end{aligned}$$

7/26

DP - EPPF

In particular, for the Dirichlet process we have that the EPPF is equal to

$$p_k^{(n)}(n_1, \dots, n_k) = \frac{a^k}{(a)_n} \prod_{j=1}^k (n_j - 1)! \quad (\alpha)_n = \frac{\Gamma(a+n)}{\Gamma(a)}$$

Does it satisfy the previous properties?

- i) $p_1^{(1)} = \frac{a^1}{(a)_1} (1-1)! = \frac{a}{a} = 1$
- ii) $p_k^{(n)}(n_1, \dots, n_k) = \frac{a^k}{(a)_n} \prod_{j=1}^k (n_j - 1)! = \frac{a^k}{(a)_n} \prod_{j=1}^k (n_{\sigma(j)} - 1)! = p_k^{(n)}(n_{\sigma(1)}, \dots, n_{\sigma(k)})$
- iii) Recall that $(a)_b = \frac{\Gamma(a+b)}{\Gamma(a)}$, then we have

$$\begin{aligned} \sum_{\ell=1}^k p_k^{(n+1)}(n_1, \dots, n_\ell + 1, \dots, n_k) + p_{k+1}^{(n+1)}(n_1, \dots, n_k, 1) \\ = \sum_{\ell=1}^k \frac{a^k}{(a)_{n+1}} \prod_{j=1}^k (n_j - 1)! \frac{(n_\ell)!}{(n_\ell - 1)!} + \frac{a^{k+1}}{(a)_{n+1}} \prod_{j=1}^k (n_j - 1)! \\ = \left(\frac{a^k \Gamma(a)}{\Gamma(a+n)} \prod_{j=1}^k (n_j - 1)! \right) \left[\frac{\Gamma(a+n)}{\Gamma(a+n+1)} \left(a + \sum_{\ell=1}^k n_j \right) \right] = p_k^{(n)}(n_1, \dots, n_k) \end{aligned}$$

8/26

DP - EPPF

Notice that we can derive the allocation probabilities of the Chinese restaurant process starting from the EPPF. Suppose we have allocated X_1, \dots, X_n into k distinct groups, with frequencies n_1, \dots, n_k .

The $n+1$ elements can be allocated to the j -th group with probability

$$\frac{p_k^{(n+1)}(n_1, \dots, n_j + 1, \dots, n_k)}{p_k^{(n)}(n_1, \dots, n_k)} = \frac{\frac{a^k}{(a)_{n+1}} \prod_{j=1}^k (n_j - 1)!}{\frac{a^k}{(a)_n} \prod_{j=1}^k (n_j - 1)!} = \frac{n_j}{a+n},$$

and to a new group with probability

$$\frac{p_{k+1}^{(n+1)}(n_1, \dots, n_k, 1)}{p_k^{(n)}(n_1, \dots, n_k)} = \frac{\frac{a^{k+1}}{(a)_{n+1}} \prod_{j=1}^k (n_j - 1)!}{\frac{a^k}{(a)_n} \prod_{j=1}^k (n_j - 1)!} = \frac{a}{a+n}.$$

The previous probability are equivalent to the allocation process described by a CRP.

9/26

Dirichlet process mixtures

Idea of mixtures of distributions:
if one model is not flexible enough
→ Use more than one model

Introduction

Let \mathbb{Y}, Θ be Polish spaces, and \mathcal{W}, \mathcal{F} their Borel σ -field. Let $K(\cdot; \cdot)$ be a (regular enough) kernel function, such that $K : \mathbb{X} \times \Theta \rightarrow [0, 1]$, $K(\cdot; \theta)$ is a density function for all $\theta \in \Theta$, and $K(A; \cdot)$ is measurable for all $A \subseteq \mathbb{Y}$. We can define a mixture model in a general fashion as

$$f(y) = \int_{\Theta} K(y; \theta) P(d\theta)$$

Examples:

- $P = \sum_{j=1}^k p_j \delta_{x_j}(\cdot)$, finite number of elements with known and fixed weights p_j 's and atoms x_j 's;
- $P = \sum_{j=1}^k W_j \delta_{X_j}(\cdot)$, finite number of elements with unknown random weights W_j 's and atoms X_j 's;
- $P = \sum_{j=1}^{\infty} W_j \delta_{X_j}(\cdot)$, infinite number of elements with unknown random weights W_j 's and atoms X_j 's;

10/26

DP mixture

We set $P \sim DP(a, P_0(\cdot))$, with $a > 0$, and $P(\cdot) = \sum_{j=1}^{\infty} W_j \delta_{\theta_j}(\cdot)$. Then we obtain an infinite mixture model

$$\tilde{f}(y) = \int_{\Theta} K(y; \theta) P(d\theta) = \sum_{j=1}^{\infty} W_j K(y; \theta_j)$$

where the last holds thanks to the a.s. discreteness of the Dirichlet process, further we know that $\{\theta_j\}_{j \geq 1} \stackrel{iid}{\sim} P_0(\cdot)$ and $\{W_j\}_{j \geq 1}$ is a sequence characterized by a stick-breaking representation $W_j = V_j \prod_{\ell < j} (1 - V_\ell)$, and $V_j \sim \text{Beta}(1, a)$ for all $j \geq 1$.

Suppose we collected a sample Y_1, \dots, Y_n distributed as the previous mixture. Then we have, in hierarchical form:

$$\begin{aligned} Y_i | \theta_i &\stackrel{\text{ind}}{\sim} K(Y_i; \theta_i) \\ \theta_1, \dots, \theta_n | P &\stackrel{\text{iid}}{\sim} P \\ P &\sim DP(a; P_0(\cdot)) \end{aligned}$$

This model is defining a random density (which is a distribution over the densities space)
It's not a fixed model, it's a random density!
Its realizations are densities themselves! The source of randomness is given by the weights and atoms.

Each Y_i is associated to a latent parameter θ_i . The sequence $(\theta_1, \dots, \theta_n)$ can possibly show ties, and its distribution is governed by a Dirichlet process.

11/26

DP mixture

Choosing $K(\cdot, \cdot)$ equal to a continuous density, a DP mixture is a random density, in the sense that its describing a random variable taking values on $S_f \subseteq \mathbb{R}_Y$, where \mathbb{R}_Y denotes the space of densities with support \mathbb{Y} .

Different choices for $K(\cdot, \cdot)$ lead to different models. Some common choices:

- $K(y; \theta) \sim N(\mu, \sigma^2)$, with $\theta = \mu$, location Gaussian mixture model.
- $K(y; \theta) \sim N(\mu, \sigma^2)$, with $\theta = (\mu, \sigma^2)$, location-scale Gaussian mixture model.
- $K(y; \theta) \sim N(\mu, \Sigma)$, with $\theta = (\mu, \Sigma)$, multivariate location-scale Gaussian mixture model.

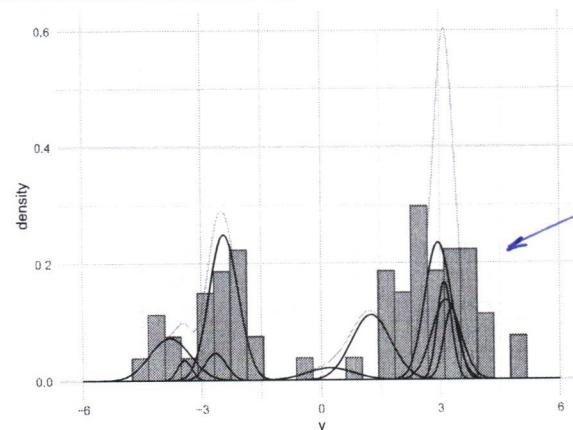
A mixture of Gaussian distribution is a really flexible model – dense in a space of regular densities (Lo 1984).

You can define a mixture model for several kind of data:

- Functional data – ex. $K(y, \theta) \sim GP(\mu(t), R_\theta(t, t'))$ with $\theta = (\mu(t), R_\theta(t, t'))$
- Graph data – ex. $K(y, \theta) \sim ERGM(\tau)$ with $\theta = \tau$ (Exponential Random Graph Model).
- ...

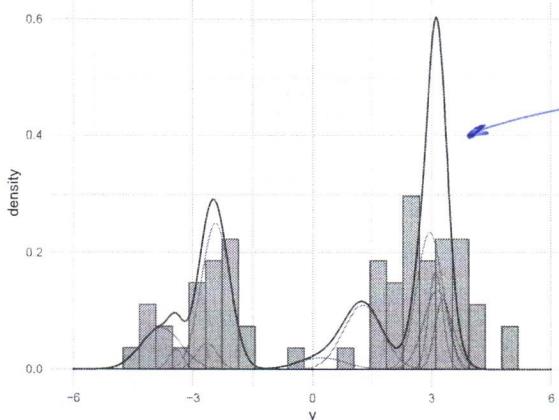
12/26

Example of location-scale Gaussian mixture



The black lines are different components in the mixture. Each observation is associated to a specific component. We have a finite number of different components with random weights, random location and random variances.

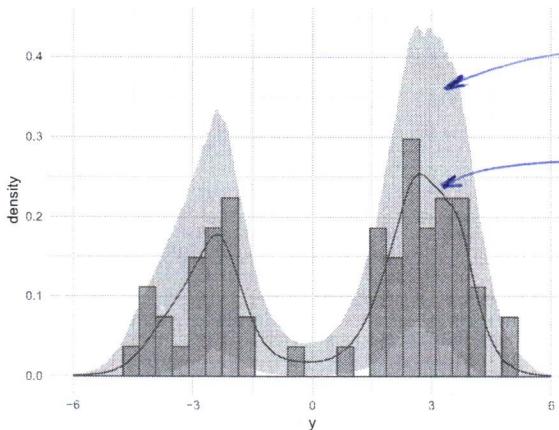
13/26

Example of location-scale Gaussian mixture

single realization from the random density that we induce with this mixture model (unweighted version of the previous components: we're taking the average of the distrib. weighted by their random weights)

→ A realization from the mixture model is a single mixture (which is a particular density)

14/26

Example of location-scale Gaussian mixture

credible bands (0.9)

posterior mean

15/26

See the GIF!**Latent partition in the mixtures**

Let Y_1, \dots, Y_n be a set of random elements taking values on \mathbb{Y} . We can augment the problem by introducing a set of latent random variables S_1, \dots, S_n , where the generic $S_i \in \mathbb{N}$ describe the allocation of the i -th observation to a specific group, i.e. $S_i = j$ if the i -th observation belongs to the j -th component of the mixture, equivalently

$$S_i = j \quad \text{if} \quad \theta_i = \theta_j^*$$

We are implicitly inducing equivalence classes in our data, such that

$$Y_i \sim Y_\ell \quad \text{if} \quad S_i = S_\ell$$

or

$$Y_i \sim Y_\ell \quad \text{if} \quad \theta_i = \theta_\ell$$

The group allocations S_1, \dots, S_n are unknown and object of interest.

We are implicitly assuming a latent random partition of the data $\Psi_n = \{A_1, \dots, A_k\}$ with $A_j = \{i : S_i = j\}$.

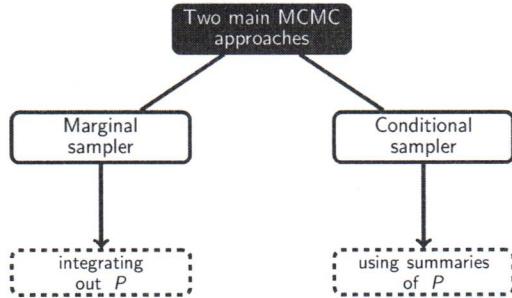
- Model based clustering → we assign the observations to the best model.
- Product partition model: $\Psi_n = \{A_1, \dots, A_k\}$ with $A_j = \{i : S_i = j\}$

$$P(\Psi_n | \mathbf{Y}, \boldsymbol{\theta}) \propto \pi(\Psi_n) P(\boldsymbol{\theta} | \Psi_n) L(\mathbf{Y} | \Psi_n, \boldsymbol{\theta})$$

standard procedure
our procedure: we're (sort of) considering ∞ clusters: when our clusters are not sufficient to describe the data we've taking a new cluster

16/26

Sampling the mixtures



17/26

Sampling the mixtures

We are interested into perform inference on the latent partition and the random density underlying the data.

Our target distribution is to sample from:

$$\begin{aligned}\mathcal{L}(\mathbf{Y}, \boldsymbol{\theta}, P) &= \mathcal{L}(\mathbf{Y} | \boldsymbol{\theta}, P) \mathcal{L}(\boldsymbol{\theta} | P) \mathcal{L}(P) \\ &= \mathcal{L}(\mathbf{Y} | \boldsymbol{\theta}) \mathcal{L}(\boldsymbol{\theta} | P) \mathcal{L}(P)\end{aligned}$$

Two possible strategies:

1. Marginal sampler, we resort to

$$E_P[\mathcal{L}(\mathbf{Y}, \boldsymbol{\theta})] = \mathcal{L}(\mathbf{Y} | \boldsymbol{\theta}) E_P[\mathcal{L}(\boldsymbol{\theta} | P)]$$

i.e. we marginalize out the distribution of P .

2. Conditional samplers, we update $P \sim \mathcal{L}(P)$ and then, conditionally on the current value of P we sample from

$$\mathcal{L}(\mathbf{Y} | \boldsymbol{\theta}) \mathcal{L}(\boldsymbol{\theta} | P)$$

Different methods – different estimation of the uncertainty associated to the problem.

18/26

Marginal sampler

We are interested into update the cluster allocation of a single observation, conditionally on the rest of the sample. Let $\mathbf{Y}_{(i)} = \{Y_1, \dots, Y_n\} \setminus Y_i$ and $\boldsymbol{\theta}_{(i)} = \{\theta_1, \dots, \theta_n\} \setminus \theta_i$. Let $\theta_{1(i)}, \dots, \theta_{k_{(i)}(i)}$ be the $k_{(i)}$ distinct values in $\boldsymbol{\theta}_{(i)}$, with frequencies $n_{1(i)}, \dots, n_{k_{(i)}(i)}$. Then we have

$$\begin{aligned}\mathcal{L}(Y_i, \theta_i | \mathbf{Y}_{(i)}, \boldsymbol{\theta}_{(i)}) &= \int_{\mathbb{P}_\Theta} \mathcal{L}(Y_i, \theta_i, P | \mathbf{Y}_{(i)}, \boldsymbol{\theta}_{(i)}) dP \\ &= \int_{\mathbb{P}_\Theta} \underbrace{\mathcal{L}(Y_i, \theta_i | P)}_{Y_i | \theta_i \stackrel{\text{ind}}{\sim} \mathcal{K}(Y_i; \theta_i), \theta_i | P \sim P} \underbrace{\mathcal{L}(P | \boldsymbol{\theta}_{(i)})}_{P \sim D^P(\alpha(\cdot) + \sum_{j=1, j \neq i}^n \delta_{\theta_j(\cdot)})} dP \\ &= \int_{\mathbb{P}_\Theta} \int_{\Theta} \mathcal{K}(Y_i; \theta_i) P(d\theta_i) \mathcal{L}(P | \boldsymbol{\theta}_{(i)}) dP \\ &= \int_{\Theta} \mathcal{K}(Y_i; \theta_i) \int_{\mathbb{P}_\Theta} P(d\theta_i) \mathcal{L}(dP | \boldsymbol{\theta}_{(i)}) \\ &= \int_{\Theta} \mathcal{K}(Y_i; \theta_i) E[P(d\theta_i) | \boldsymbol{\theta}_{(i)}] \\ &= \frac{a}{a+n} \int_{\Theta} \mathcal{K}(Y_i; \theta_i) P_0(\theta_i) + \frac{n}{a+n} \sum_{j=1}^{k_{(i)}} \frac{n_{j(i)}}{a+n} \mathcal{K}(Y_i; \theta_{j(i)}^*)\end{aligned}$$

probability of the i -th observation of being assigned to a new group

Marginal sampler

The marginal sampler is based on

$$\mathcal{L}(Y_i, \theta_i | \mathbf{Y}_{(i)}, \boldsymbol{\theta}_{(i)}) = \frac{a}{a+n} \int_{\Theta} \mathcal{K}(Y_i; \theta_i) P_0(\theta_i) + \frac{n}{a+n} \sum_{j=1}^{k_{(i)}} \frac{n_{j(i)}}{a+n} \mathcal{K}(Y_i; \theta_{j(i)}^*)$$

probability of the i -th observation of being assigned to an already existing group

Remarks:

- ▶ Choose $P_0(\cdot)$ in a smart way! If $P_0(\cdot)$ is conjugate to $\mathcal{K}(\cdot, \cdot)$, then the integral $\int_{\Theta} \mathcal{K}(Y_i; \theta_i) P_0(\theta_i)$ is the prior predictive distribution, often available in closed form. Otherwise there are other strategies, for example you can evaluate it via Monte Carlo integration.
- ▶ When you re-allocate the i -th observation, and it is a singleton, its cluster must be destroyed (this is the reason for the nefarious notation in the previous slide).
- ▶ The distribution, which correspond to the predictive distribution of a DPM, is still a convex combination between the prior part and an empirical part, weighted by the kernel.
- ▶ We can trace the group allocation using the variables S_1, \dots, S_n with

$$P(S_i = j | \text{rest}) \propto n_j K(Y_i; \theta - j(i))$$

$$P(S_i = k_{(i)} + 1 | \text{rest}) \propto a \int_{\Theta} K(Y_i; \theta) P_0(d\theta)$$

20/26

Marginal sampler

Algorithm 1: Pseudocode for the marginal sampler

```

1 set admissible initial values  $\theta^{(0)}$ ;
2 for each iteration  $r = 1, \dots, R$  do
3   for each  $i = 1, \dots, n$  do
4     let  $k_{(i)}$  be the number of distinct values  $\theta_{j(i)}^{*(r)}$ 's in
         $\theta_{(i)}^{(r)} = (\theta_1^{(r)}, \dots, \theta_{j-1}^{(r)}, \theta_{j+1}^{(r)}, \dots, \theta_n^{(r)})$  and  $n_{j(i)}$ , for  $j = 1, \dots, k_{(i)}$ , the
        corresponding frequencies;
5     sample  $S_i^{(r)}$  from
         $P(S_i^{(r)} = t | -) \propto \begin{cases} n_j \mathcal{K}(X_i; \theta_{j(i)}^{*(r)}) & \text{if } j \in \{1, \dots, k_{(i)}\} \\ a \int \mathcal{K}(X_i; \theta) P_0(d\theta) & \text{if } j = k_{(i)} + 1 \end{cases}$ 
6     for each unique value  $\theta_j^{*(r)}$  in  $\theta^{(r)}$  do
7       update  $\theta_j^{*(r)}$  from  $P(\theta_j^{*(r)} \in dt | -) \propto P_0(dt) \prod_{i: S_i=j} \mathcal{K}(X_i; t)$ 
8 end

```

Another approach:

21/26

Truncated stick-breaking sampler

We are interested in both update the cluster allocation all single observation, and estimate the random density. Then in a Gibbs style algorithm, we have

$$\begin{aligned} \mathcal{L}(\mathbf{Y}, \mathbf{S}, P) &= \mathcal{L}(\mathbf{Y} | \mathbf{S}, P) \mathcal{L}(\mathbf{S} | P) \mathcal{L}(P) \\ &= \mathcal{L}(\mathbf{Y} | \mathbf{S}) \mathcal{L}(\mathbf{S}) \mathcal{L}(P) \end{aligned}$$

where

$$\begin{aligned} \mathcal{L}(\mathbf{Y} | \mathbf{S}, P) &= \prod_{i=1}^n \mathcal{K}(Y_i; \theta_i) = \prod_{j=1}^k i \prod_{i: S_i=j} \mathcal{K}(Y_i; \theta_j^*), \\ \mathcal{L}(\mathbf{S}) &\propto \prod_{i=1}^n \sum_{j=1}^M \mathbf{1}_{[S_i=j]} \end{aligned}$$

and P is distributed as a truncated DP

$$P = \sum_{j=1}^M W_j \delta_{\theta_j^*}(\cdot)$$

with $W_j = V_j \prod_{\ell < j} (1 - V_\ell)$, $V_j \sim \text{Beta}(1, a)$ if $j < M$, $V_M = 1$ and $\{\theta_1^*, \dots, \theta_M^*\} \stackrel{\text{iid}}{\sim} P_0(\cdot)$.

22/26

Truncated stick-breaking sampler

Conditionally on \mathbf{Y}, θ, P we can reallocate the observations, where

$$P(S_i = j | -) \propto W_j \mathcal{K}(Y_i; \theta_j^*),$$

we can then update P as

$$\mathcal{L}(P | \mathbf{Y}, \mathbf{S}) = \sum_{j=1}^M W_j \delta_{\theta_j^*}(\cdot)$$

with

$$W_j = V_j \prod_{\ell < j} (1 - V_\ell),$$

and

$$V_j \sim \text{Beta}(1 + n_j, a + m_j) \text{ if } j < M, \quad V_M = 1$$

where $n_j = \sum_{i=1}^n \mathbf{1}_{[S_i=j]}$ and $m_j = \sum_{i=1}^n \mathbf{1}_{[S_i>j]}$. We can update the θ_j^* 's from

$$\mathcal{L}(\theta_j^* | \mathbf{Y}, \mathbf{S}) \propto P_0(d\theta_j^*) \prod_{i: S_i=j} \mathcal{K}(Y_i; \theta_j^*)$$

23/26

Truncated stick-breaking sampler

Algorithm 2: Pseudocode for the truncated stick-breaking sampler

```

1 set admissible initial values  $\mathbf{S}^{(0)}, P_M$ ;
2 for each iteration  $r = 1, \dots, R$  do
3   for each  $i = 1, \dots, n$  do
4     sample  $S_i^{(r)}$  from
         $P(S_i^{(r)} = j | -) \propto W_j^{(r)} \mathcal{K}(X_i; \theta_{j(i)}^{*(r)}), \quad j = 1, \dots, M$ 
5     for  $j = 1, \dots, M$  do
6       update  $V_j$  from  $\mathcal{L}(V_j^{(r)} | -) \sim \text{Beta}(1 + n_j, a + m_j)$  if  $j < M$ ,  $V_M = 1$ 
7       with  $n_j = \sum_{i=1}^n \mathbf{1}_{[S_i=j]}$  and  $m_j = \sum_{i=1}^n \mathbf{1}_{[S_i>j]}$ 
8       set  $W_j^{(r)} = V_j^{(r)} \prod_{\ell < j} (1 - V_\ell^{(r)})$  from
9       update  $\theta_j^{*(r)}$  from
         $\mathcal{L}(\theta_j^{*(r)} | \mathbf{Y}, \mathbf{S}) \propto P_0(d\theta_j^{*(r)}) \prod_{i: S_i=j} \mathcal{K}(Y_i; \theta_j^{*(r)})$ 
10 end

```

24/26

To synthesize

Different sampling strategies to deal with P

- **Marginal sampler:** we integrate out the distribution of P , resorting to the distribution of $\theta_i \mid \theta_1, \dots, \theta_{i-1}, \theta_{i+1}, \dots, \theta_n$.
- **Truncated stick-breaking sampler:** we work with finite summaries of P .

Note that:

- With the **marginal sampler** we focus on $\mathcal{L}(\theta_1, \dots, \theta_n) = E[\mathcal{L}(\theta_1, \dots, \theta_n \mid P)]$, the update must be done sequentially.
- With the **conditional samplers** we focus on $\mathcal{L}(\theta_1, \dots, \theta_n \mid P)$, the update can be done in parallel.

Both the algorithms produce a realization for the latent partition of the data at each iteration → we need a point estimate.

25/26

Example

Mixture of Gaussian distribution.

We consider $K(Y; \theta) \sim N(\mu, \sigma^2)$, with $\theta = \mu$, i.e. location mixtures.

We choose as base measure a Gaussian distribution, $P_0(\mu) \sim N(m_0, \sigma_0^2)$, to preserve conjugacy, i.e.

$$\int K(y; \theta) P_0(d\theta) \sim N(m_0, \sigma_0^2 + \sigma^2)$$

We consider a model

$$\tilde{f}(y) = \sum_{j=1}^{\infty} W_j K(y; \theta_j^*)$$

We compare the marginal sampling strategy and the truncated stick-breaking sampling strategy.

26/26

See the code!

10/12/2020

5/6

```
# R & STAN are friends!
library(coda)
library(rBnFmix)
library(purrr)
library(eggc1)
require(gplots)
require(eggpor)

# for plots
library(ggplot2)
library(tidy)
library(dplyr)

# MARGINAL SAMPLER
MS_univ <- function(data, niter, burnin, m0, sigma2_0, sigma2, mass){
  n <- length(data)
  res_part <- matrix(0, ncol = n, nrow = niter - burnin)
  mu <- m0

  # initialize the partition - all in one cluster
  S <- rep(1, n)
  mu_out <- list()
  probs <- c()

  #-----#
  # allocation step
  #-----#
  pb <- txtProgressBar(0, niter, style = 3)
  # for all the iterations
  for(iter in 1:niter){

    #-----#
    # for all the observations
    for(i in 1:n){

      # J = number of current clusters
      J <- length(mu)
      probs <- rep(0, J + 1)

      # for all the active clusters
      for(j in 1:J){
        prob[i,j] <- (sum(S[-j] == j) / (n + mass - 1)) * dnorm(data[i], mean = mu[j], sd = sqrt(sigma2))
      }

      # probability of new cluster
      prob[i, J + 1] <- (mass / (n + mass - 1)) * dnorm(data[i], mean = m0, sd = sqrt(sigma2 + sigma2_0))

      # sample from 1:J+1 with probability prob
      S[i] <- sample(1:(J + 1), size = 1, replace = F, prob = probs)

      # if the cluster is new, generate its mean value
      if(S[i] > J){
        mu[S[i]] <- rnorm(1, mean = (data[i] / sigma2 + m0 / sigma2_0) * (1 / sigma2 + 1 / sigma2_0),
                           sd = sqrt(1 / (1 / sigma2 + 1 / sigma2_0)))
      }

      # resize the objects
      mu <- mu[unique(sort(S))]
      S <- as.numeric(as.factor(S))

      if(iter > burnin){
        res_part[iter - burnin, ] <- S
        mu_out[iter - burnin] <- mu
      }

      #-----#
      # acceleration step
      #-----#
    }
  }

  #-----#
  # for all the objects
  for(j in 1:length(mu)){
    ntemp <- length(data[S == j])
    mu[j] <- rnorm(1, mean = (sum(data[S == j]) / sigma2 + m0 / sigma2_0) * (length(data[S == j]) / sigma2 + 1 / sigma2_0),
                  sd = sqrt(1 / (length(data[S == j]) / (1 / sigma2 + 1 / sigma2_0))))
  }

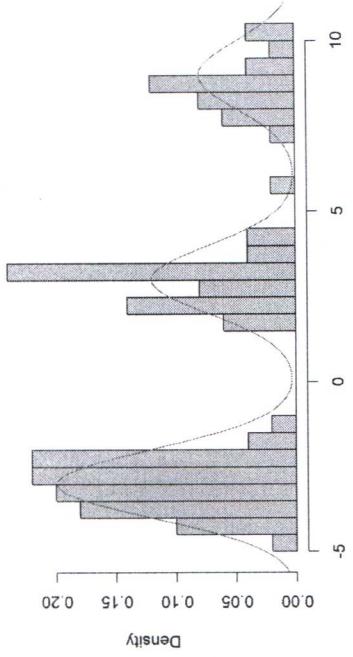
  setTxtProgressBar(pb, iter)
}

close(pb)
return(list(res_part, mu_out))
}

# TRUNCATED SB SAMPLER
TRSB_univ <- function(data, niter, burnin, m0, sigma2_0, sigma2, mass, N){

  n <- length(data)
```

Histogram of data



```

res_part <- matrix(0, ncol = n, nrow = niter - burnin)
mu <- me0
w <- v <- c()

# Initialize the partition - all in one cluster
S <- rep(1, n)

mu_out <- list()
probs_out <- list()
probs <- c()

pb <- txtProgressBar(0, niter, style = 3)
# for all the iterations
for(iter in 1:niter){
  #-----
  # UPDATE THE WEIGHTS and THE AIMS
  #-----

  for(j in 1:M){
    ntemp <- length(data[S == j])
    mu[j] <- rnorm(1, mean = (sum(data[S == j]) / sigma2_0) + (length(data[S == j]) / sd
      mass2 + 1 / sigma2_0), sd = sqrt(1 / (length(data[S == j]) / sigma2_0 + 1 / sigma2_0)))
  }

  V[1] <- rbeta(1, 1 + sum(S == 1), mass + sum(S > 1))
  w[1] <- V[1]

  for(i in 2:(M - 1)){
    V[i] <- rbeta(1, 1 + sum(S == 1), mass + sum(S > i))
    w[i] <- V[i] * prod(1 - V[1:(i - 1)])
  }

  V[M] <- 1
  w[M] <- V[M] * prod(1 - V[1:(M - 1)])
  #-----

  # allocation step
  #-----
}

# for all the observations
for(i in 1:n){

  # for all the active clusters
  for(j in 1:M){

    probs[j] <- w[j] * dnorm(data[i], mean = mu[j], sd = sqrt(sigma2))

  }
}

# sample from 1:M with probability probs
S[1] <- sample(1:M, size = 1, replace = F, prob = probs)
}

# if(iter > burnin){
  res_part[iter - burnin, ] <- S
  mu_out[[iter - burnin]] <- mu
  probs_out[[iter - burnin]] <- w
}

return(list(res_part, mu_out))
}

#-----
#-----#
# with simulated data
data <- c(rnorm(50, -3, 1), rnorm(30, 3, 1), rnorm(20, 9, 1))
hist(data, breaks = 25, freq = F)
sq <- seq(-6, 13, by = 0.1)
fsq <- .5 * dnorm(sq, -3, 1) + .3 * dnorm(sq, 3, 1) + .2 * dnorm(sq, 9, 1)
lines(sq, fsq, col = 2)

```

10/12/2020

6/6

```
## Fraction in 1st window = 0.1
## Fraction in 2nd window = 0.5
## CL_MAR CL_TRSB
## 0.1344 -1.7143

acfplot(coda_import)

# WHY NOT WITH THE PARAMETERS?
#-----#
#-----#
#-----#
# a first visualization of the clustering
# compute the posterior similarity matrix (PSM)
# then plot the heatmap
PSM <- function(cls){
  n <- ncol(cls)
  mat_out <- matrix(0, ncol = n, nrow = n)
  for(i in 1:n){
    for(j in 1:i){
      mat_out[i,j] <- mat_out[j,i] <- mean(cls[,i] == cls[,j])
    }
  }
  return(mat_out)
}

heatmap(PSM(sample_MAR[[1]]), Rowv = NA, Colv = NA)
heatmap(PSM(sample_TRSB[[1]]), Rowv = NA, Colv = NA)

dev.off()

## null device
## 1

par(mfrow = c(3,1))
plot(1:1000, sample_TRSB[[1]][,91], type = "l")
plot(1:1000, sample_TRSB[[1]][,51], type = "l", col = 2)
plot(1:1000, sample_TRSB[[1]][,1], type = "l", col = 3)
dev.off()

## null device
## 1

# diagnostic on the number of clusters
CL_MAR <- apply(sample_MAR[[1]], 1, function(x) length(unique(x)))
CL_TRSB <- apply(sample_TRSB[[1]], 1, function(x) length(unique(x)))

coda_import <- as.mcmc(cbind(CL_MAR, CL_TRSB))

summary(coda_import)

## Iterations = 1:1000
## Thinning interval = 1
## Number of chains = 1
## Sample size per chain = 1000
##
## 1. Empirical mean and standard deviation for each variable,
## plus standard error of the mean:
##
##          Mean   SD Naive SE Time-series SE
## CL_MAR  6.395 1.667  0.05270   0.06342
## CL_TRSB 4.783 1.639  0.05183   0.04134
##
## 2. Quantiles for each variable:
##
##          2.5% 25% 50% 75% 97.5%
## CL_MAR   4     5     6     7     10
## CL_TRSB  2     4     5     6     8

Genekie.diag(coda_import)
```

BNP3

Riccardo Corradin
December 11, 2020

1/15

Recap

Let $\mathbb{N}_n = \{1, \dots, n\}$ space of first n integers.

Definition

We say that $\psi_n = \{A_1, \dots, A_k\}$ is a *partition* of \mathbb{N}_n if the sets $\{A_1, \dots, A_k\}$ are nonempty and such that $\mathbb{N}_n = \cup_i A_i$ and $A_i \cap A_j = \emptyset$ for all $i \neq j$.

Let \mathcal{B} be the space of all partitions of \mathbb{N}_n and \mathcal{B} its discrete σ -field.

Definition

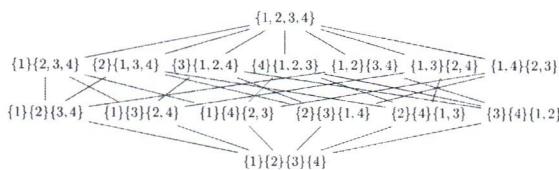
Let (Ω, \mathcal{F}, P) be a generic probability space. A random partition is any measurable function $\Psi_n : (\Omega, \mathcal{F}, P) \rightarrow (\mathcal{B}, \mathcal{B})$.

The definition of Ψ_n induces a probability measure on the space of partitions \mathcal{B} .

2/15

Recap

Problem: the space of the possible partition of n elements is large, and is growing fast as far as n grows.



We have

$$S_{n,k} = \frac{1}{k!} \sum_{j=0}^k (-1)^j \binom{k}{j} (k-j)^n$$

(Stirling number of the second kind) ways to partition n elements in k groups, and

$$B_n = \sum_{k=1}^n S_{n,k}$$

(Bell number) ways to partition n elements. B_n is exploding as far as n grows.

3/15

Recap

We are in a mixture model framework. Suppose we have a set $\mathbf{Y} = \{Y_1, \dots, Y_n\}$ of observations, and we assume a mixture model

$$\tilde{f}(y) = \int_{\Theta} \mathcal{K}(y; \theta) P(d\theta)$$

with $P \sim DP(a, P_0(\cdot))$. Equivalently

$$\begin{aligned} Y_i | \theta_i &\stackrel{\text{ind}}{\sim} \mathcal{K}(Y_i; \theta_i) \\ \theta_1, \dots, \theta_n | P &\stackrel{\text{iid}}{\sim} P \\ P &\sim DP(a, P_0(\cdot)) \end{aligned}$$

P a.s. discrete \rightarrow the sequence $\theta_1, \dots, \theta_n$ may show ties
 \rightarrow latent partition Ψ_n in $\theta_1, \dots, \theta_n$

$\pi(\Psi_n) = p_k^{(n)}(n_1, \dots, n_k)$ EPPF of a Dirichlet process.

We saw how to sample from the posterior distribution of a DP mixture.

$$\tilde{f}(Y | Y_1, \dots, Y_n) = \int_{\Theta} \mathcal{K}(Y; \theta) P(d\theta | \theta_1, \dots, \theta_n)$$

At each step we have a partition of the data, i.e. a realization from $\pi(\Psi_n | \mathbf{Y})$.

4/15

Recap

Scan the entire partitions space is unfeasible in a reasonable time, as far as $n \nearrow$.

n	1	2	3	4	5	6	7	8	9	10
B_n	1	2	5	15	52	203	877	4140	21147	115975

Given a set of observations, we aim to estimate the best partition latent in the data among the B_n possible partitions.

We set a prior $\pi(\Psi_n)$ on the partitions' space.

We are interested in area where the posterior concentrates its mass.

We need a strategy to reduce the number of partitions considered in the estimation, starting from our MCMC realizations.

5/15

Point estimate

We approach the problem from a decision theory point of view. Let $L(\cdot, \cdot)$ denotes a loss function, two partitions of n elements as arguments. We have that the optimal partition $\psi_n^* \in \mathbb{B}$ (or equivalently $(S_1, \dots, S_n) \in \mathbb{N}^n$) is the solution of

$$\begin{aligned}\psi_n^* &= \arg \min_{\psi_n \in \mathbb{B}} \left\{ \mathbb{E}_{\psi_n} [L(\psi_n, \hat{\psi}_n) \mid \mathbf{Y}] \right\} \\ &= \arg \min_{\psi_n \in \mathbb{B}} \left\{ \sum_{\psi_n \in \mathbb{B}} L(\psi_n, \hat{\psi}_n) \mathcal{L}(\Psi_n = \psi_n \mid \mathbf{Y}) \right\}\end{aligned}$$

where ψ_n^* denotes the optimal partition, $\hat{\psi}_n$ the partition we are considering, ψ_n the partitions wrt we take the expectation, and $\mathcal{L}(\psi_n \mid \mathbf{Y})$ is the posterior probability of ψ_n , given a set of observations \mathbf{Y} .

- The space \mathbb{B} is too large
- We can choose $L(\cdot, \cdot)$ in several ways
- We can set different priors on Ψ_n , for example the EPPF of a DP

$L(\cdot, \cdot)$ = how much we're losing describing one partition in terms of the other (it's in terms of the distance between the two partitions)

6/15

1. Point estimate - 0 – 1 = dummy loss function

A first loss function that we can consider is the 0 – 1 loss function, i.e.

$$L_{0-1}(\psi_n, \hat{\psi}_n) = \mathbf{1}_{[\psi_n = \hat{\psi}_n]}$$

Such loss function is taking value 1 when all the blocks of ψ_n and $\hat{\psi}_n$ coincide, and 0 otherwise.

- Is not accommodating possible similarities in the partitions, all the partitions different from $\hat{\psi}_n$ are penalized in the same way
- the point estimate coincides with

$$\arg \min_{\psi_n \in \mathbb{B}} \left\{ \sum_{\psi_n \in \mathbb{B}} \mathbf{1}_{[\psi_n = \hat{\psi}_n]} \mathcal{L}(\Psi_n = \psi_n \mid \mathbf{Y}) \right\} = \arg \max_{\psi_n \in \mathbb{B}} \left\{ \mathcal{L}(\Psi_n = \hat{\psi}_n \mid \mathbf{Y}) \right\}$$

the maximum a posteriori of the distribution of Ψ_n .

7/15

2. Point estimate - Binder

A more relaxed loss function: the Binder loss function.

$$L_B(\psi_n, \hat{\psi}_n) = \sum_{j < i} \left[C_1 \mathbf{1}_{[S_i = S_j]} \mathbf{1}_{[S_i \neq \hat{S}_j]} + C_2 \mathbf{1}_{[S_i \neq S_j]} \mathbf{1}_{[S_i = \hat{S}_j]} \right]$$

and by setting $C_1 = C_2$ it can be written as

$$L_B(\psi_n, \hat{\psi}_n) = \frac{1}{2} \left(\sum_{i=1}^{k_n} n_{i*}^2 + \sum_{j=1}^{\hat{k}_n} n_j^2 - 2 \sum_{j=1}^{\hat{k}_n} \sum_{i=1}^{k_n} n_{ij}^2 \right)$$

it is based on the cluster frequencies, where

- k_n is the number of blocks in ψ_n
- \hat{k}_n is the number of blocks in $\hat{\psi}_n$
- $n_{ij} = \sum_{\ell=1}^n \mathbf{1}_{[S_\ell=i]} \mathbf{1}_{[\hat{S}_\ell=j]}$ is the number of observations in the i -th block of ψ_n and in the j -th block of $\hat{\psi}_n$
- $n_{i*} = \sum_{j=1}^{\hat{k}_n} n_{ij} = \sum_{\ell=1}^n \mathbf{1}_{[S_\ell=i]}$ is the number of observations in the i -th block of ψ_n
- $n_{*j} = \sum_{i=1}^{k_n} n_{ij} = \sum_{\ell=1}^n \mathbf{1}_{[\hat{S}_\ell=j]}$ is the number of observations in the j -th block of $\hat{\psi}_n$

here we're penalizing with C_1 if two elements agree in ψ_n and disagree in $\hat{\psi}_n$ and we penalize with C_2 if two elements agree in $\hat{\psi}_n$ but disagree in ψ_n .

8/15

Point estimate - Binder

You can show that the solution of

$$\psi_n^* = \arg \min_{\hat{\psi}_n \in \mathbb{B}} \left\{ \sum_{\psi_n \in \mathbb{B}} L_B(\psi_n, \hat{\psi}_n) \mathcal{L}(\psi_n = \hat{\psi}_n | \mathbf{Y}) \right\}$$

is given by the partition which, in binary representation, minimize the distance from the posterior similarity matrix \mathcal{M} , where

$$\mathcal{M}_{ij} = P(S_i = S_j | \mathbf{Y})$$

and

$$\mathbf{S}^* = \arg \min_{\mathbf{S} \in \mathbb{N}^n} \left\{ \sum_{i < j} |\mathbf{1}_{[\hat{S}_i = \hat{S}_j]} - \mathcal{M}_{ij}| \right\}$$

The Binder loss function, differently from the $0-1$ loss function, is not penalizing in the same way the partitions different from the optimal one.

9/15

Point estimate - Binder

A practical example: consider $n = 3$, and all the possible partitions are

$$\mathbb{B} = \{\{1, 2, 3\}, \{1, 2\}, \{1, 3\}, \{2, 3\}, \{1\}, \{2\}, \{3\}\}$$

Equivalently in a binary representation

$$A_1 = \begin{bmatrix} 1 \\ 0 & 1 \\ 0 & 0 & 1 \end{bmatrix}, \quad A_2 = \begin{bmatrix} 1 \\ 1 & 1 \\ 0 & 0 & 1 \end{bmatrix}, \quad A_3 = \begin{bmatrix} 1 \\ 0 & 1 \\ 0 & 1 & 1 \end{bmatrix},$$

$$A_4 = \begin{bmatrix} 1 \\ 0 & 1 \\ 1 & 0 & 1 \end{bmatrix}, \quad A_5 = \begin{bmatrix} 1 \\ 1 & 1 \\ 1 & 1 & 1 \end{bmatrix}$$

posterior prob. that the observation 1 and 2 are grouped together

$$\text{and } \mathcal{M} = \begin{bmatrix} 1 \\ 0.8 & 1 \\ 0.2 & 0.2 & 1 \end{bmatrix}, \text{ then we have}$$

$$d(A_1, \mathcal{M}) = 0.8 + 0.2 + 0.2 = 1.2; \quad d(A_2, \mathcal{M}) = 0.6; \quad d(A_3, \mathcal{M}) = 1.8;$$

$$d(A_4, \mathcal{M}) = 1.8; \quad d(A_5, \mathcal{M}) = 2.4$$

and the point estimate is given by A_2 .

10/15

3. Point estimate - Variation of information

Another loss function: variation of information.

$$L_{VI}(\psi_n, \hat{\psi}_n) = H(\psi_n) + H(\hat{\psi}_n) - 2I(\psi_n, \hat{\psi}_n)$$

where $H(\cdot)$ denotes the entropy and $I(\cdot, \cdot)$ the relative entropy, and it can be written as

$$L_{VI}(\psi_n, \hat{\psi}_n) = \sum_{i=1}^{k_n} \frac{n_i}{n} \log \left(\frac{n_i}{n} \right) + \sum_{j=1}^{\hat{k}_n} \frac{n_j}{n} \log \left(\frac{n_j}{n} \right) - 2 \sum_{j=1}^{\hat{k}_n} \sum_{i=1}^{k_n} \frac{n_{ij}}{n} \log \left(\frac{n_{ij}}{n} \right)$$

it is based on the cluster frequencies, where

- k_n is the number of blocks in ψ_n
- \hat{k}_n is the number of blocks in $\hat{\psi}_n$
- $n_{ij} = \sum_{\ell=1}^n \mathbf{1}_{[S_\ell=i]} \mathbf{1}_{[\hat{S}_\ell=j]}$ is the number of observations in the i -th block of ψ_n and in the j -th block of $\hat{\psi}_n$
- $n_i = \sum_{j=1}^{\hat{k}_n} n_{ij} = \sum_{\ell=1}^n \mathbf{1}_{[S_\ell=i]}$ is the number of observations in the i -th block of ψ_n
- $n_j = \sum_{i=1}^{k_n} n_{ij} = \sum_{\ell=1}^n \mathbf{1}_{[\hat{S}_\ell=j]}$ is the number of observations in the j -th block of $\hat{\psi}_n$

11/15

Point estimate

We still have a problem: the partitions' space is too large.

We restrict our attention to a subset of the entire partition space $\mathbb{S}_R \subseteq \mathbb{B}$, the space of the partition visited during R steps of an MCMC algorithm.

Our solution is sub-optimal, but feasible

$$\psi_n^* = \arg \min_{\hat{\psi}_n \in \mathbb{B}} \left\{ \mathbb{E}_{\psi_n} [L(\psi_n, \hat{\psi}_n) | \mathbf{Y}] \right\} \simeq \arg \min_{\hat{\psi}_n \in \mathbb{S}_R} \left\{ \mathbb{E}_{\psi_n} [L(\psi_n, \hat{\psi}_n) | \mathbf{Y}] \right\}$$

Instead of considering B_n elements, we consider only R partitions. You can also reduce this number due to the fact that some partitions are visited more than one time.

- This strategy can be relaxed by scanning the neighborhood of the partitions visited by the sampler, with a *greedy* strategy;
- You can also provide an estimation of the uncertainty associated to the point estimate by defining a credible ball on the partition space.

Instead of considering the entire partition space we restrict our attention just to the subset of the partitions that we visit during the MCMC estimation. We saw that at each iteration we produce a partition: instead of considering the whole partition space we consider the R partitions that we visited during the estimation of the model

12/15

Examples

Location Gaussian DP mixture:

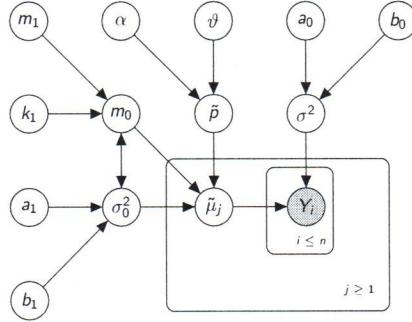
$$\begin{aligned} \mathcal{K}(Y; \mu, \sigma^2) &\stackrel{\text{ind}}{\sim} N(\mu, \sigma^2) \\ \underline{\mu \mid P \stackrel{\text{iid}}{\sim} P} \\ P &\sim DP(a, P_0) \\ P_0 &\sim N(m_0, \sigma_0^2) \\ \sigma^2 &\sim IG(a_0, b_0) \\ (m_0, \sigma_0^2) &\sim NIG(m_1, k_1, a_1, b_1) \end{aligned}$$

Location-scale Gaussian DP mixture:

$$\begin{aligned} \mathcal{K}(Y; \mu, \sigma^2) &\stackrel{\text{ind}}{\sim} N(\mu, \sigma^2) \\ \underline{(\mu, \sigma^2) \mid P \stackrel{\text{iid}}{\sim} P} \\ P &\sim DP(a, P_0) \\ P_0 &\sim NIG(m_0, k_0, a_0, b_0) \\ m_0 &\sim N(m_1, \sigma_1^2) \\ k_0 &\sim G(\tau_1, \zeta_1) \\ b_0 &\sim G(a_1, b_1) \end{aligned}$$

13/15

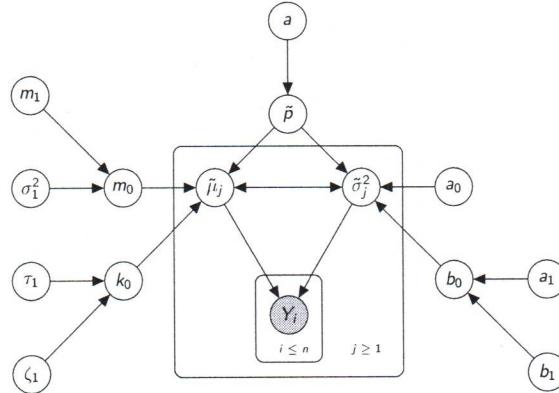
Examples



Hierarchical representation of the univariate location DP mixture model with Gaussian kernel and Gaussian base measure.

14/15

Examples



Hierarchical representation of the univariate location-scale DP mixture model with Gaussian kernel and Normal-Inverse-Gamma

15/15

See the code!

```

# R & STAN are friends!
library(coda)
library(BHmix)

# for plots
library(ggplot2)
library(tidyverse)
library(dplyr)

require(ggpubr)

library(purrr)
library(ggsci)
require(gplots)

# MARGINAL SAMPLER
MS_univ <- function(data, niter, burnin, m0, sigma2_0, sigma2, mass){

  n <- length(data)

  res_part <- matrix(0, ncol = n, nrow = niter - burnin)
  mu <- m0

  # initialize the partition - all in one cluster
  S <- rep(1, n)

  mu_out <- list()
  probs <- c()

  pb <- txtProgressBar(0, niter, style = 3)

  # for all the iterations
  for(iter in 1:niter){

    #-----#
    # allocation step
    #-----#

    # for all the observations
    for(i in 1:n){

      # J = number of current clusters
      j <- length(mu)
      probs <- rep(0, j + 1)

      # for all the active clusters
      for(j in 1:j){
        probs[j] <- (sum(S[-i] == j) / (n + mass - 1)) * dnorm(data[i], mean = mu[,j], sd = sqrt(sigma2))
      }

      # probability of new cluster
      probs[j + 1] <- (mass / (n + mass - 1)) * dnorm(data[i], mean = m0, sd = sqrt(sigma2 + sigma2_0))

      # sample from 1:j with probability probs
      S[i] <- sample(1:(j + 1), size = 1, replace = F, prob = probs)

      # if the cluster is new, generate its mean value
      if(S[i] > j){
        mu[S[i]] <- rnorm(1, mean = (data[i] / sigma2 + m0 / sigma2_0) * (1 / sigma2 + 1 / sigma2_0))
        S[i] <- sort(1 / (1 / sigma2 + 1 / sigma2_0))
      }
    }

    # resize the objects
    mu <- mu[unique(sort(S))]
    S <- as.numeric(as.factor(S))

    if(iter > burnin){
      res_part$burnin[j] <- S
      mu_out[[iter - burnin]] <- mu
    }

    #-----#
    # acceleration step
    #-----#
  }

  for(j in 1:length(mu)){
    ntemp <- length(data[S == j])
    mu[j] <- rnorm(1, mean = (sum(data[S == j]) / sigma2 + m0 / sigma2_0) * (length(data[S == j]) / sigma2 + 1 / sigma2_0),
                  + 1 / sigma2_0),
    sd = sqrt(length(data[S == j]) / (1 / sigma2 + 1 / sigma2_0))
  }

  setTxtProgressBar(pb, iter)
}

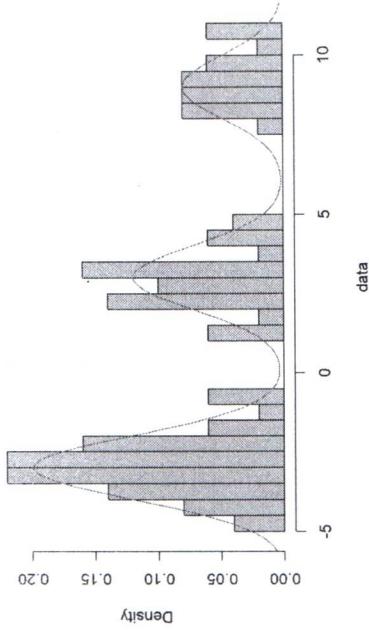
# TRUNCATED SB SAMPLER
TRSB_univ <- function(data, niter, burnin, m0, sigma2_0, sigma2, mass, M){

  n <- length(data)
  res_part <- matrix(0, ncol = n, nrow = niter - burnin)
  mu <- m0
  W <- V <- c()

}

```

Histogram of data



```

# initialize the partition - all in one cluster
S <- rep(1, n)

mu_out <- list()
probs_out <- list()
probs <- c()

pb <- txtProgressBar(0, niter, style = 3)
for(j in 1:M){
  # for all the iterations
  for(ier in 1:niter){
    #-----#
    # UPDATE THE WEIGHTS and THE ATOMS
    #-----#
    for(j in 1:M){
      temp <- length(data[S == j])
      mu[j] <- rnorm(1, mean = (sum(data[S == j]) / sigma2 + mu / sigma2_theta) * (length(data[S == j]) / sigma2))
      + 1 / sigma2_theta),
      sd = sqrt(1 / (length(data[S == j]) / sigma2 + 1 / sigma2_theta)))
    }

    V[1] <- rbeta(1, 1 + sum(S == 1), mass + sum(S > 1))
    W[1] <- V[1]

    for(l in 2:(M - 1)){
      V[l] <- rbeta(1, 1 + sum(S == l), mass + sum(S > l))
      W[l] <- V[l] * prod(1 - V[1:(l - 1)])
    }

    V[M] <- 1
    W[M] <- V[M] * prod(1 - V[1:(M - 1)])
  }
  #-----#
  # allocation step
  #-----#
}

# for all the observations
for(i in 1:n){
  # sample from i:4 with probability probs
  S[i] <- sample(1:M, size = 1, replace = F, prob = probs)
}

# for all the active clusters
for(j in 1:M{
  probs[j] <- W[j] * dnorm(data[i], mean = mu[j], sd = sqrt(sigma2))
}

#-----#
# with simulated data
data <- c(rnorm(59 - 3, 1), rnorm(38, 3, 1), rnorm(28, 9, 1))
hist(data, breaks = 25, freq = F)
sq <- seq(-6, 13, by = .5 * sqrt(.5 * dnorm(-4, 0, 1) + .3 * dnorm(sq, 3, 1) + .2 * dnorm(sq, 9, 1)))
lines(sq, fsq, col = 2)
}
#-----#
#-----#

```

```
## Fraction in 1st window = 0.1
## Fraction in 2nd window = 0.5
## CL_MAR CL_TRSB
## 2.2153 -0.1163
```

```
acfplot(coda_import)
```

```
# WHY NOT WITH THE PARAMETERS?
```

```
#-----
```

```
#-----
```

```
#-----
```

```
# a first visualization of the clustering
```

```
# compute the posterior similarity matrix (PSM)
```

```
# then plot the heatmap
```

```
PSM <- function(cls){
  n <- ncol(cls)
  mat_out <- matrix(0, ncol = n, nrow = n)
  for(i in 1:n){
    for(j in 1:i){
      mat_out[i,j] <- mat_out[j,i] <- mean(cls[,i] == cls[,j])
    }
  }
  return(mat_out)
}
```

```
heatmap(PSM(sample_MAR[[1]]), Rowv = NA, Colv = NA)
```

```
heatmap(PSM(sample_TRSB[[1]]), Rowv = NA, Colv = NA)
```

```
#-----
```

```
#-----
```

```
#-----
```

```
est_part_BINDER <- function(clust, PSM){
```

```
res_loss <- c()
```

```
for(i in 1:nrow(clust)){
```

```
binary_part <- sapply(clust[,i], function(x) as.numeric(x == clust[,i]))
```

```
res_loss[i] <- sum(abs(binary_part - PSM))
```

```
}
```

```
return(which.min(res_loss))
```

}

```
part <- matrix(c(
```

```
1,2,1,1,2,3,
```

```
1,1,1,2,2,2,
```

```
1,2,1,2,2,3,
```

```
3,1,3,3,1,1,
```

```
3,1,3,3,2,2,
```

```
), ncol = 6, byrow = T)
```

```
PSM(part)
```

```
#-----
```

```
#-----
```

```
#-----
```

```
#-----
```

```
#-----
```

```
#-----
```

```
#-----
```

```
#-----
```

```
#-----
```

```
#-----
```

```
#-----
```

```
#-----
```

```
#-----
```

```
#-----
```

```
#-----
```

```
#-----
```

```
#-----
```

```
#-----
```

```
#-----
```

```
#-----
```

```
#-----
```

```
#-----
```

```
#-----
```

```
#-----
```

```
#-----
```

```
#-----
```

```
#-----
```

```
#-----
```

```
#-----
```



```
## null device
## 1
par(mfrow = c(3,1))
plot(1:2500, sample_TRSB[[1]][,91], type = "l")
plot(1:2500, sample_TRSB[[1]][,51], type = "l", col = 2)
plot(1:2500, sample_TRSB[[1]][,1], type = "l", col = 3)
dev.off()
```

```
## null device
## 1
```

```
# diagnostic on the number of clusters
CL_MAR <- apply(sample_MAR[[1]], 1, function(x) length(unique(x)))
CL_TRSB <- apply(sample_TRSB[[1]], 1, function(x) length(unique(x)))
coda_import <- as.mcmc(cbind(CL_MAR, CL_TRSB))
summary(coda_import)
```

```
## Iterations = 1:25000
## Thinning interval = 1
## Number of chains = 1
## Sample size per chain = 25000
```

```
## 1. Empirical mean and standard deviation for each variable,
```

```
## plus standard error of the mean:
```

```
##          Mean   SD Naive SE Time-series SE
## CL_MAR 6.417 1.555 0.03169 0.03310
## CL_TRSB 4.793 1.485 0.02970 0.02436
##          [1,] 6.417 1.555 0.03169 0.03310
##          [2,] 4.793 1.485 0.02970 0.02436
##          [3,] 6.417 1.555 0.03169 0.03310
##          [4,] 4.793 1.485 0.02970 0.02436
##          [5,] 6.417 1.555 0.03169 0.03310
##          [6,] 4.793 1.485 0.02970 0.02436
```

```
geweke.diag(coda_import)
```

```
est_part_BINDER(part, PSM(part))
```

```
[1] 1 2 1 1 2 3
```


11/12/2020

6/16

```

## Iterations = 1:2500
## Thinning interval = 1
## Number of chains = 1
## Sample size per chain = 2500
## Empirical mean and standard deviation for each variable,
## plus standard error of the mean:
##      Mean     SD   Naive SE Time-series SE
## 1       1       0       0       0
## 2. Quantiles for each variable:
##    2.5%  25%  50%  75%  97.5%
## 1       1       1       1       1
## 2       1       1       1       1       1

plot(coda_MAR_mv)
effectivesize(coda_MAR_mv)

## var1
## 0

summary(model_MAR_mv)

## Pdensity function call:
## 2500 burn-in iterations
## 5000 iterations
## Global estimation time: 18.23 seconds
## Average number of groups: 1
## Min number of groups: 1 ; Max number of groups: 1
## 2. Quantiles for each variable:
##    2.5%  25%  50%  75%  97.5%
## 3       3       4       4       5

# plots
# p1 <- plot(model_MAR_mv, show_points = T, show_clust = T, dimension = c(1,2))
# p2 <- plot(model_MAR_mv, show_points = T, show_clust = T, dimension = c(1,3))
# p3 <- plot(model_MAR_mv, show_points = T, show_clust = T, dimension = c(2,3))
# ggarrange(p1, p2, NULL, p3, ncol = 2, nrow = 2)
part_MAR_mv <- partition(model_MAR_mv, dist = "VI")
table(part_MAR_mv$partitions[,1])

## 
## 1
## 139

part_MAR_mvH <- partition(model_MAR_mv, dist = "Binder")
table(part_MAR_mv$partitions[,1])

## 
## 1
## 139

mcmc = list(nitriter = 5000, nburn = 2500,
            method = "MCMC", model = "LS", hyper = T, print_message = T),
          output = list(grid = expand_grid(seq(-4, 4, by = 0.2),
                                          seq(-4, 4, by = 0.2),
                                          seq(-4, 4, by = 0.2)), out_type = "MEAN"))

model1_MAR_mvH <- Pdensity(Clust_data,
                                mcmc = mcmc, nburn = 5000, nitriter = 5000, nburn = 2500,
                                method = "MCMC", model = "LS", hyper = T, print_message = T),
                                output = list(grid = expand_grid(seq(-4, 4, by = 0.2),
                                          seq(-4, 4, by = 0.2),
                                          seq(-4, 4, by = 0.2)), out_type = "MEAN"))

## 
## 1
## 17
## 10

part_MAR_mvH <- partition(model_MAR_mvH, dist = "Binder")
table(part_MAR_mv$partitions[,1])

## 
## 1
## 2
## 17
## 10

```