# BAYESIAN INFERENCE

We start knowing
- $\theta$ (or $\Theta$): the parameters about the population characteristics that we want to estimate
- $Y$: the numerical description of the sample that we get (to interpret)

So we have $\theta(\Theta) \in \Theta$, the parameter space, and $Y \in Y$, the sample space.

We start as follows:

(1) we assign a prior distribution about how we believe initially that the population characteristics behave

$$\theta \sim \pi(\theta)$$

(2) we assign the likelihood $f(Y|\theta)$, ie the conditional law of a sample obs $Y$ given the true parameter $\theta$

$$Y_1, \ldots, Y_m | \theta \overset{iid}{\sim} f_{Y|\theta}(Y|\theta)$$

$$\Rightarrow f_{Y|\theta}(Y|\theta) = f_{Y|\theta}(Y_1, \ldots, Y_m|\theta) = \prod_{i=1}^{m} f_{Y_i|\theta}(Y_i|\theta)$$

(3) once we collect the data $Y$ we can get the posterior distribution which updates the belief about $\theta$. We can compute it by the Bayes thm:

$$\pi(\theta|Y) = \frac{f(Y|\theta)\cdot\pi(\theta)}{f(Y)} = \left|_{\text{argument track}} \frac{f(Y|\theta)\cdot\pi(\theta)}{\int_\Theta f(Y|\theta)\cdot\pi(\theta)\,d\theta}\right. =$$

marginal density of $Y$
which we also call $m(Y)$

$$= \left|_{\text{iid}} \frac{\prod_i f(Y_i|\theta)\cdot\pi(\theta)}{\int_\Theta \prod_i f(Y_i|\theta)\pi(\theta)\,d\theta}\right. \propto (\text{LIKELIHOOD})\,(\text{PRIOR})$$

for the computation we can just rule out the den, to leave the functional part, the num

A very interesting example was about the posterior computation of $\theta$ when we have/observe two samples $Y_1$ and $Y_2$. But morally, however we proceed we get the same distribution in the end.

Case of using both $Y_1$ and $Y_2$ together

$$\pi(\theta|Y_1, Y_2) = \frac{f(Y_1, Y_2|\theta)\cdot\pi(\theta)}{f(Y_1, Y_2)} = \left|_{\text{cond}} \frac{f(Y_1|\theta)\cdot f(Y_2|\theta)\cdot\pi(\theta)}{f(Y_1, Y_2)}\right. =$$

$$= \frac{f(Y_2|\theta)}{f(Y_2|Y_1)} \cdot \underbrace{\frac{f(Y_1|\theta)\cdot\pi(\theta)}{f(Y_1)}}_{\pi(\theta|Y_1)} = \frac{f(Y_2|\theta)\cdot\pi(\theta|Y_1)}{f(Y_2|Y_1)} =$$

$$= \pi(\theta|Y_1, Y_2) \qquad \text{Case of first using } Y_1 \text{ and then } Y_2$$

## INFERENTIAL PROBLEMS

(1) Bayesian point estimation. we could take the posterior distribution (always refer to the posterior when dealing uncertainty) and we could take the posterior mean or median, for example.
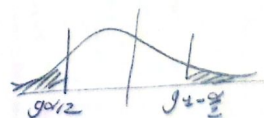
$$\hat\theta = \arg\min_\theta E[\ell(\theta, \hat\theta)|Y] \qquad \ell(\theta, \hat\theta) = \begin{cases} (\theta-\hat\theta)^2 \\ |\theta-\hat\theta| \\ 1-\delta_\theta(\hat\theta) \end{cases} \Rightarrow \hat\theta = \begin{cases} E[\theta|Y] \\ \text{median}(\theta|Y) \\ \text{MAP}(\theta|Y) \end{cases}$$

(2) Bayesian interval estimation. Here we have two methods, where the target was to build ORC $\theta$ st $P(\theta \in CR|Y) \geq 1-\alpha$, with $\alpha$ low. we call it CR as credible region (or CI as credible interval).

• we could take the quantiles to delimit a $1-\alpha$ probability lot)
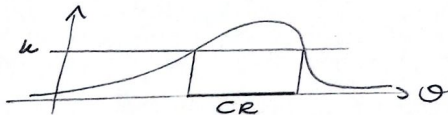
$$CR_{\theta|Y} = \left(g_{\frac{\alpha}{2}}^{\pi(\theta|Y)}, g_{1-\frac{\alpha}{2}}^{\pi(\theta|Y)}\right) \quad \text{wf } \theta \in CR$$

$$CR_{\theta;|Y} = \left(g_{\frac{\alpha}{2}}^{\pi;(\theta|Y)}, g_{1-\frac{\alpha}{2}}^{\pi;(\theta|Y)}\right) \quad \text{wf } \theta \in CR^{p_0}$$



$g_{\alpha/2}$     $g_{1-\frac{\alpha}{2}}$

15

- we could take the HPD region (highest posterior density region) which would be the narrowest region/interval that meets that request. For a uni-modal distr the rules make sense:

$$CR = \{\theta \in \Theta : \pi(\theta | \underline{x}) \geq k\}$$
$$k : P(\theta \in CR | \underline{x}) = 1 - \alpha$$



(3) Hypothesis testing. We could have two cases, the comparison set comparison or the parameter target.

- $\begin{cases} H_0 | \theta \in \Theta_0 \\ H_1 | \theta \in \Theta_1 \end{cases}$  prior: on $\pi(\theta)$ or the more precise  on $\begin{cases} f_0(\theta) | & \text{if } \theta \in \Theta_0 \\ f_1(\theta) | & \text{if } \theta \in \Theta_1 \end{cases}$

Let's call $\pi_0 = \pi(\Theta_0)$ and $\pi_1 = \pi(\Theta_1) = 1 - \pi_0$ since $\Theta_1 = \Theta_0^c$. These will give the prior odds. For the posterior we surely need the likelihood / marginal

$$f(\underline{x}) = \int_\Theta f(\underline{x}|\theta)\pi(\theta) \, d\theta =$$
$$= \pi_0 \int_{\Theta_0} f(\underline{x}|\theta)f_0(\theta) \, d\theta + \pi_1 \int_{\Theta_1} f(\underline{x}|\theta)f_1(\theta) \, d\theta$$

$$\Rightarrow \pi(\theta|\underline{x}) = \begin{cases} \dfrac{\pi_0 f(\underline{x}|\theta)f_0(\theta)}{f(\underline{x})} & \text{if } \theta \in \Theta_0 \\[3mm] \dfrac{\pi_1 f(\underline{x}|\theta)f_1(\theta)}{f(\underline{x})} & \text{if } \theta \in \Theta_1 \end{cases}$$

Now we could just use 'ea example if $P(\theta \in \Theta_0 | \underline{x}) \geq^? 0.5$ or similar, but the rules we ras compare also with the prior class.

$$(\text{POST ODDS}) = \frac{P(\theta \in \Theta_0 | \underline{x})}{P(\theta \in \Theta_1 | \underline{x})} = \frac{\pi_0}{\pi_1} \cdot \frac{\int_{\Theta_0} f(\underline{x}|\theta)f_0(\theta) \, d\theta}{\int_{\Theta_1} f(\underline{x}|\theta)f_1(\theta) \, d\theta} =$$

$$= \frac{\pi_0}{\pi_1} \cdot BF_{01} = (\text{PRIOR ODDS}) \cdot BF_{01}$$

- $\begin{cases} H_0 | \theta = \theta_0 \\ H_1 | \theta \neq \theta_0 \end{cases}$  for this case the rules is the same so the above where we can set as for the following

$$\theta \sim \begin{cases} f_0(\theta) = \delta_{\theta_0}(\theta) & \text{if } \theta = \theta_0 \\ f_1(\theta) = \pi(\theta) & \text{if } \theta \neq \theta_0 \end{cases} \Rightarrow BF_{01} = \frac{f(\underline{x}|\theta_0)}{f(\underline{x})}$$

Then once we have $BF_{01}$ we can interpret wrt according to this table from literature:

$$-\log_{10}(BF_{01}) \in \begin{cases} (0, 1/2) & \text{barely mentionable} \\ (1/2, 1) & \text{substantial} \\ (1, 2) & \text{strong} \\ (2, +\infty) & \text{decisive} \end{cases} \quad \left(\begin{array}{c}\text{evidence} \\ \text{for } H_0\end{array}\right)$$

(4) Posterior predictive distribution (while the prior pred distr would just be the marginal, we $f(\underline{y}) = \int_\Theta f(\underline{y}|\theta)\pi(\theta) \, d\theta$). Here we want to get the law of a new obs given all the others. So

$$f_{Y_{m+1} | Y_1 \dots Y_m}(\underline{y}|\underline{x}) = \frac{f(\underline{y}, \underline{x})}{f(\underline{x})} = \bigg|_{\substack{\text{augment} \\ \text{truth}}} \frac{\int_\Theta f(\underline{y}, \underline{x}|\theta)\pi(\theta) \, d\theta}{\int_\Theta f(\underline{x}|\theta)\pi(\theta) \, d\theta} =$$

$$= \bigg|_{\text{cond}} \int_\Theta f(\underline{y}|\theta) \left( \frac{f(\underline{x}|\theta)\pi(\theta)}{\int_\Theta f(\underline{x}|\theta)\pi(\theta) \, d\theta} \right) d\theta = \int_\Theta f(\underline{y}|\theta)\pi(\theta|\underline{x}) \, d\theta$$

In another way, gucha:

$$f_{Y_{m+1}|\underline{x}}(\underline{y}|\underline{x}) = \bigg|_{\substack{\text{augment} \\ \text{truth}}} \int_\Theta f_{Y_{m+1}|\underline{x}}(\underline{y}, \theta | \underline{x}) \, d\theta =$$

$$= \int_\Theta f(\underline{y}|\theta, \underline{x}) \cdot \pi(\theta|\underline{x}) = \bigg|_{\text{cond}} \int_\Theta f(\underline{y}|\theta) \cdot \pi(\theta|\underline{x}) \, d\theta$$

# PRIORS STUFF

We say that the r.v.s $(y_1, \ldots, y_m)$ are exchangeable wf for any permutation $\pi$ of $(1, \ldots, m)$, we have that $\mathcal{L}(y_1, \ldots, y_m) = \mathcal{L}(y_{\pi(1)}, \ldots, y_{\pi(m)})$. This condition implies that the order in which data are recorded is irrelevant for inferential purposes.

This implies a (symmetry) about the role of the individuals in the sampling, like a "homogeneous conditions". It is relevant for Bayesian stuff through a

**Thm** (De Finetti's representation)

a $(0,1)$ sequence $(y_m)_{m \geq 1}$ of binary r.v.s is exchangeable

$(\Leftrightarrow)$ $\exists$ a prob measure $F$ on $([0,1], \mathcal{B}([0,1]))$ st

$$p_{\underline{y}}(\underline{y}) = \int_0^1 \left( \theta^{\Sigma y_i} (1-\theta)^{m - \Sigma y_i} \right) F(\theta) \, d\theta$$

$\forall m \geq 1$ and $\forall \underline{y} \in \{0,1\}$

$(\Leftrightarrow)$ $\exists$ some r.v. $\tilde{\theta}$ st we can model

$$y_1, \ldots, y_m \mid \tilde{\theta} \overset{iid}{\sim} Ber(\tilde{\theta})$$
$$\tilde{\theta} \sim F(\theta)$$

About actually defining priors we can have different choices.
(1) **Reference** (or convenience) priors. We use them when we want to have a minimal impact on the Bayesian analysis. They are also called "non-informative", but this definition is misleading, they are more a reliable/common choice to start model with.

An example is a flat prior, $\pi(\theta) = c \; \forall \theta \in \Theta$. This works well if $\Theta$ is bounded, otherwise becomes improper. But a flat prior with a representation may not be flat in another representation. $\sim (\theta \sim Beta(1,1)$ and $\gamma = \ln(\theta/1-\theta))$

(2) **Jeffreys** priors. They are often improper, but are invariant for monotone transformations. They are defined through the Fisher information

$$I(\theta) = E\left[ \left( \frac{\partial}{\partial \theta} \ln f(y|\theta) \right)^2 \right] = E\left[ -\frac{\partial^2}{\partial \theta^2} \ln f(y|\theta) \right] \quad \substack{E \, wrt \, y, \\ not \, \theta}$$

$$\pi(\theta) \propto \sqrt{I(\theta)} \quad \sim \substack{\text{from the proportional} \\ \text{we take just the functional} \\ \text{part of } \theta, \text{ discard constants}}$$
$$\rightarrow \text{we can use both } f(y|\theta) \text{ and } f(\underline{y}|\theta) \text{ to compute wt}$$

(3) **Informative** (or scientifically informed) priors. The best choice to use, especially for delicate/relevant params of the model.

In any way, we could also be more free in this choice as eventually, having more data, all the priors will never on ways (converge) to the result posterior. We

**Thm** (AN of the posterior distr). Let $y_1, \ldots, y_m \mid \theta \overset{iid}{\sim} f(y|\theta)$ and $\pi(\theta)$ the prior of $\theta \in \Theta \subseteq \mathbb{R}^d$. Under suitable regularity conditions we have that

$$\pi(\theta|\underline{y}) \underset{m \to +\infty}{\approx} N(\tilde{\theta}_m, V_m) \qquad \substack{\tilde{\theta}_m: \text{ posterior mean} \\ V_m: \text{ posterior cov matrix}}$$

# POPULAR DISTRIBUTIONS

$X \sim Pois(\lambda) \;(\Leftrightarrow)\; p_X(u) = e^{-\lambda} \frac{\lambda^u}{u!} \mathbb{1}_{\mathbb{N}_0}(u) \qquad \substack{E(x) = \lambda \\ var(x) = \lambda}$

$X \sim Exp(\lambda) \;(\Leftrightarrow)\; f_X(x) = \lambda e^{-\lambda x} \mathbb{1}_{(0,+\infty)}(x) \qquad \substack{E(x) = 1/\lambda \\ var(x) = 1/\lambda^2}$

$X \sim Beta(a,b) \quad \Longleftrightarrow \quad f_X(x) = \frac{\Gamma(a+b)}{\Gamma(a)\Gamma(b)} x^{a-1}(1-x)^{b-1} \mathbb{1}_{[0,1]}(x)$

Comments:

- $X_\omega \overset{\perp}{\sim} \Gamma(a_\omega, b) \implies \frac{X_\omega}{X_1+X_2} \sim Beta(a_1, a_2)$
- It's also an extension of the $U([0,1])$ law. To show model different moments make clear the first last one

$E(X) = \frac{a}{a+b}$

$var(X) = \frac{a \cdot b}{(a+b+1)(a+b)^2}$

$S^{2-1} = S^1 = \{x \in \mathbb{R}: x \in [0,1], 0 \leq x \leq 1\}$

---

$\begin{pmatrix} x_1 \\ \vdots \\ x_m \end{pmatrix} \sim Dir\left(\alpha = \begin{pmatrix} a_1 \\ \vdots \\ a_m \end{pmatrix}\right) \quad \Longleftrightarrow \quad f_{\underline{X}}(\underline{x}) = \frac{\Gamma(\Sigma a_\omega)}{\prod \Gamma(a_\omega)} \prod_{\omega=1}^{m} x_\omega^{a_\omega-1} \mathbb{1}_{S^{m-1}}\left(\underline{x} = \begin{pmatrix} x_1 \\ \vdots \\ x_{m-1} \end{pmatrix}\right)$

$\boxed{\left(\prod_{i=1}^{m-1} x_\omega^{a_\omega-1}\right)\left(1-\sum_{\omega=1}^{m-1} x_\omega\right)^{a_m-1}}$

the last component is determined by $x_1,...,x_{m-1}$

Comments:

- The support is $S^{m-1}(\underline{x}) = \{(x_1,...,x_{m-1}) \in \mathbb{R}^{m-1}: x_\omega \in [0,1], 0 \leq \Sigma x_\omega \leq 1\}$

and then $x_m$ will be $1 - \Sigma_{\omega=1}^{m-1} x_\omega$

$E(X_\omega) = \frac{a_\omega}{a_0}$

$var(X_\omega) = \frac{a_\omega(a_0 - a_i)}{(a_0+1) a_0^2}$

$(a_0 = \Sigma a_\omega)$

- It's the extension of the Beta law

- $U_\omega \overset{\perp}{\sim} \Gamma(a_\omega, \beta) \quad \Longleftrightarrow \quad X_\omega = \frac{U_\omega}{\Sigma U_\omega}, \quad \underline{X} \sim Dir(\underline{\alpha})$ for $\omega=1,...,m$

- $\begin{pmatrix} x_1 \\ x_2 \\ x_3 \end{pmatrix} \sim Dir\begin{pmatrix} a_1 \\ a_2 \\ a_3 \end{pmatrix} \implies \begin{pmatrix} x_1+x_3 \\ x_2 \end{pmatrix} \sim Dir\begin{pmatrix} a_1+a_3 \\ a_2 \end{pmatrix}$ and any

- $\underline{X} \sim Dir(\underline{\alpha}) \implies \underbrace{x_1}_{\sim Beta(a_1, a_0-a_1)} \perp \underbrace{\left(\frac{x_2}{1-x_1}, ..., \frac{x_m}{1-x_1}\right)}_{\sim Dir(a_2,...,a_m)}$

---

$X \sim Gamma(\alpha, \lambda) \quad \Longleftrightarrow \quad f_X(x) = \frac{\lambda^\alpha}{\Gamma(\alpha)} x^{\alpha-1} e^{-\lambda x} \mathbb{1}_{[0,+\infty)}(x)$

Comments:

- $\Gamma(\alpha+1) = \alpha\Gamma(\alpha)$ and $\Gamma(m+1) = m!$
- $X_\omega \overset{\perp}{\sim} \Gamma(a_\omega, \lambda) \implies \Sigma X_\omega \sim \Gamma(\Sigma a_\omega, \lambda)$ for $\omega=1,...,m$
- $\Gamma(1/2) = \sqrt{\pi}$
- $c \cdot \Gamma(\alpha, \lambda) \sim \Gamma(\alpha, \lambda/c)$ for $c>0$

$E(X) = \alpha/\lambda$
$var(X) = \alpha/\lambda^2$

$\boxed{\begin{array}{c} Y = \frac{1}{X} \sim InvGamma(\alpha, \lambda) \\ E(Y) = \frac{\lambda}{\alpha-1} \\ \frac{\lambda^\alpha}{\Gamma(\alpha)}\left(\frac{1}{x}\right)^{\alpha+1} e^{-\lambda(\frac{1}{x})} \end{array}}$

---

$X \sim Categorical(\underline{\varphi}) \quad \Longleftrightarrow \quad \varphi_X(\omega) = p_\omega \mathbb{1}_{\{1,...,u\}}(\omega)$
$= \prod_{j=1}^{k} p_j \mathbb{1}_{\{\omega=j\}}$

$\boxed{\Gamma(a) = \int_0^{+\infty} x^{a-1} e^{-x} dx}$

Comments:

- The $supp(X) = \{1,...,u\}$, we X gets a label from 1 to $u$
- The vector $\underline{p}$ is $\underline{p} = (\varphi_1,...,\varphi_u) \in S^u$, we $\Sigma \varphi_\omega = 1$ and $\varphi_\omega \in [0,1]$ too

---

$\begin{pmatrix} x_1 \\ \vdots \\ x_u \end{pmatrix} \sim Multinomial\left(m, \begin{pmatrix} \varphi_1 \\ \vdots \\ \varphi_u \end{pmatrix}\right) \quad \Longleftrightarrow \quad \varphi_{\underline{X}}(\underline{x}) = \frac{m!}{x_1! \cdots x_u!} \prod_{\omega=1}^{u} p_\omega^{x_\omega}$

Comments:

- It is the extension of the binomial, since each of the counts the # of occurrence (successes) of extracting label $i$ (out of labels 1 to $u$) when doing $m$ experiments

- In the $supp(\underline{x})$ is the set of values $x_i \in \mathbb{N}_0$ that sum up to $m$, $\Sigma x_\omega = m$ (as we get $m$ outcomes)

$E(X_i) = m p_i$
$var(X_i) = m p_i(1-p_i)$

# FAMOUS MODELS

**Beta-Binomial model.** Useful when we work on data which are obtained through a series of Bernoulli experiments.

$$\text{prior: } \theta \sim Beta(a,b)$$
$$\text{likelih: } y_1, \ldots, y_m \mid \theta \sim Ber(\theta) \; (\Rightarrow \; y \mid \theta \sim Bin(m,\theta))$$
$$\Rightarrow \text{post: } \theta \mid \underline{y} \sim Beta\left(a + \Sigma y_i, \; b + (m - \Sigma y_i)\right)$$

$$f(\theta \mid \underline{y}) = \frac{f(\theta, \underline{y})}{f(\underline{y})} \propto f(\theta, \underline{y}) = f(\underline{y}, \theta) = f(\underline{y} \mid \theta) \cdot f(\theta) =$$

$$= \left[ \theta^{\Sigma y_i} (1-\theta)^{m - \Sigma y_i} \right] \left[ \frac{1}{B(a,b)} \theta^{a-1}(1-\theta)^{b-1} \mathbb{1}_{[0,1]}(\theta) \right] \propto$$

$$\propto \theta^{(a + \Sigma y_i) - 1} (1-\theta)^{(b + m - \Sigma y_i) - 1} \mathbb{1}_{[0,1]}(\theta) \Rightarrow \text{kernel of the unnormalized Beta}$$

**Normal-Normal model.** Useful to model data which follows a normal law of known variance, and we want to model $y$ also with a normal law.

$$\text{prior: } y \sim N(y_0, \tau^2 = var)$$
$$\text{likelih: } y_1, \ldots, y_m \mid y \sim N(y, \sigma^2 = var)$$
$$\Rightarrow \text{post: } y \mid \underline{y} \sim N(y_m, \tau_m^2)$$

$$y_m = \frac{\sigma^2 y_0 + m \tau^2 \bar{y}}{\sigma^2 + m \tau^2}, \quad \tau_m^2 = \frac{\sigma^2 \tau^2}{\sigma^2 + m \tau^2}$$

*all the likelihood are actually written $\sim$ iid ..., our data are conditionally $\perp$*

$$f(y \mid \underline{y}) \propto f(\underline{y} \mid y) \cdot f(y) = \left[ \prod_{i=1}^{m} \frac{1}{\sqrt{2\pi\sigma^2}} e^{-\frac{1}{2\sigma^2}(y_i - y)^2} \right] \cdot \left[ \frac{1}{\sqrt{2\pi\tau^2}} e^{-\frac{1}{2\tau^2}(y - y_0)^2} \right] \propto$$

$$\propto e^{-\frac{1}{2\sigma^2} \Sigma (y_i - y)^2} \; e^{-\frac{1}{2\tau^2}(y - y_0)^2} = \cdots =$$

$$= e^{-\frac{1}{2}\left[ \frac{m}{\sigma^2}(y - \bar{y})^2 + \frac{1}{\tau^2}(y - y_0)^2 \right]} = \cdots$$

**Gamma-Poisson model.** Useful when our data describes rare events that we want to count (so a Poisson law is needed) suited for them.

$$\text{prior: } \theta \sim Gamma(\alpha, \beta)$$
$$\text{likelih: } y_1, \ldots, y_m \mid \theta \sim Poi(\theta)$$
$$\Rightarrow \text{post: } \theta \mid \underline{y} \sim Gamma(\alpha + \Sigma y_i, \; \beta + m)$$

$$f(\theta \mid \underline{y}) \propto f(\underline{y} \mid \theta) \cdot f(\theta) = \left[ \prod_{i=1}^{m} e^{-\theta} \frac{\theta^{y_i}}{y_i!} \right] \left[ \frac{\beta^\alpha}{\Gamma(\alpha)} \theta^{\alpha-1} e^{-\beta\theta} \mathbb{1}_{[0,+\infty)}(\theta) \right] \propto$$

$$\propto e^{-m\theta} \theta^{\Sigma y_i} \theta^{\alpha-1} e^{-\beta\theta} \mathbb{1}_{[0,+\infty)}(\theta) =$$

$$= \theta^{(\alpha + \Sigma y_i) - 1} e^{-(\beta + m)\theta} \mathbb{1}_{[0,+\infty)}(\theta) \Rightarrow \text{kernel of the unnormalized Gamma}$$

**Dirichlet-multinomial model.** Useful when we have data values that belong to different categories/classes (on each of our data we put a Multinomial law), of which the probabilities are modeled through a Dirichlet law.

$$\text{prior: } \underline{p} \sim Dir(\underline{a}), \; \underline{p} \in \mathbb{R}^{u = \# classes}$$
$$\text{likelih: } \underline{y} = (y_1, \ldots, y_u) \mid \underline{p} \sim Mult(m, \underline{p})$$
$$\Rightarrow \text{post: } \underline{p} \sim Dir(\underline{a} + \underline{y})$$

$$\text{remember that } \underline{y} \mid \underline{p} \sim Mult(m, \underline{p}) \text{ means that}$$
$$\underline{y} = \begin{pmatrix} y_1 \\ \vdots \\ y_u \end{pmatrix}, \; y_i = \# events \text{ that belong to category } i$$
$$(m = \Sigma y_i = \# \text{ total observations})$$

$$f(\underline{p} \mid \underline{x}) \propto f(\underline{x} \mid \underline{p}) \cdot f(\underline{p}) = \left[ \frac{n!}{x_1! \dots x_u!} \prod_{i=1}^{u} p_i^{x_i} \right] \cdot \left[ \frac{1}{B(\alpha)} \prod_{i=1}^{u} p_i^{\alpha_i - 1} \, \mathcal{U}_{S_{u-1}}(\underline{p}) \right]$$

$$\propto \prod_{i=1}^{u} p_i^{(\alpha_i + x_i) - 1} \, \mathcal{U}_{S_{u-1}}(\underline{p}) \quad \Rightarrow \quad \text{kernel of the} \atop \text{updated Dir}$$

## normal, Inverse Gamma – normal model. likelihood model normal data, but now with both parameters unknown.



prior: $Y \mid \sigma^2 \sim N(Y_0, \sigma^2/\lambda)$
$\sigma^2 \sim \text{Inv Gamma}(\alpha, \beta)$ $\Rightarrow$ $Y \not\perp \sigma^2$

likel: $Y_1, \dots, Y_n \mid (Y, \sigma^2) \sim N(Y, \sigma^2)$

$\Rightarrow$ post: $Y \mid \sigma^2 \sim N\left( \frac{n\bar{Y} + \lambda Y_0}{n + \lambda}, \frac{\sigma^2}{n + \lambda} \right)$ $\Rightarrow$ $Y \not\perp \sigma^2$

$\sigma^2 \sim \text{Inv Gamma}\left( \alpha + \frac{n}{2}, \, \beta + \frac{1}{2} \sum (Y_i - \bar{Y})^2 + \frac{1}{2} \frac{n\lambda(\bar{Y} - Y_0)^2}{n + \lambda} \right)$

# SIMULATION METHODS (MCMC)

We now want to get relevant features about our target distri. $\pi(\cdot \mid \Sigma)$, the posterior, through simulations, as often we can't directly sample out from it (we would mean get values that follow a certain distribution).

Idea: resort on two thms about simulations and long-runs properties.

**Thm** (SLLN). Let $\theta^{(1)}, \theta^{(2)}, \dots$ be an iid sequence of samples from $\pi(\theta)$, and let $q: \Theta \to \mathbb{R}$ st $E[|q(\theta)|] < +\infty$. Our target is computing

$$\bar{q} = E_\pi(q(\theta)) = \int_\Theta q(\theta)\, \pi(\theta)\, d\theta$$

Then we have that

(1) $\quad \bar{q}^{(T)} = \frac{1}{T} \sum_{t=1}^{T} q(\theta^{(t)}) \xrightarrow[T\to+\infty]{a.s.} \bar{q}$

(2) $\quad y_\alpha^{(T)} = \binom{\alpha\text{-quantile}}{\text{of } q(\theta^{(t)})} \xrightarrow[T\to+\infty]{a.s.} y_\alpha = \binom{\alpha\text{-quantile}}{\text{of } q(\theta)}$

**Thm** (CLT). If we have the above assumptions, plus $\text{var}[q(\theta)] < +\infty$, then we get new info about the velocity of convergence:

set $\sigma^2 = \text{var}[q(\theta)]$
$\sigma^{2(T)} = \text{var}[q(\theta^\omega)]$ over $\omega = 1, \dots, T$

(1) $\quad \sqrt{T}\left(\bar{q}^{(T)} - \bar{q}\right) \xrightarrow[T\to+\infty]{\mathcal{D}} \mathcal{N}(0, \sigma^2) \quad (=) \quad \bar{q}^{(T)} \underset{T\to+\infty}{\approx} \mathcal{N}\left(\bar{q}, \frac{\sigma^2}{T}\right)$

(2) $\quad \sigma^{2(T)} \xrightarrow[T\to+\infty]{a.s.} \sigma^2$

## MARKOV CHAINS MONTE CARLO

The idea is to get samples from the posterior $\pi(\theta \mid \Sigma)$ by evolving a MC whose limiting distr. is that posterior, and so the samples will be given by the strategy we are writing.

Let $E \subset \mathbb{R}^u$ be the state space. A time-homogeneous MC $(X_m)_{m\geq 0}$ with values in $E$ is a sequence of r.v.s $E$-valued st

$$P(X_{m+1} \in A \mid X_m = x_m, \dots, X_0 = x_0) = P(X_{m+1} \in A \mid X_m = x_m)$$

we call $\quad P(x_m, A) = P(X_{m+1} \in A \mid X_m = x_m)$
$\quad P^m(x, A) = P(X_m \in A \mid X_0 = x)$
$\quad P_x(\cdot) = P(\cdot \mid X_0 = x)$
$\left.\right\}$ $P(\cdot, \cdot)$ is the transition prob-ability kernel

We now use more terminology and properties, and then focus on the needed requirements for the running algos.

(1) **Invariant** (or stationary) distribution (as we need also a distr. for $X_0$ to define a MC). A measure $\pi$ on $(E, \Sigma)$ st

$$\int_E \pi(x)\, P(x, A)\, dx = \pi(A) \qquad \forall A \in \Sigma \qquad (\pi P = \pi)$$
$\qquad \qquad \underline{\boxed{1}} \quad \boxed{1}$

(2) **Irreducibility**: a MC is irreducible wf $\omega_i, j$ (w_j can be reached from $\omega_i$) $\forall i, j \in E$ in a finite # of steps. This also means that eventually (sooner or later) we will visit all the reachable states / sets of states.

$$\exists \psi \text{ a meas distr}: \forall A \text{ st } \psi(A) > 0 \quad \exists m_{x,A} \geq 1: P^m(x, A) > 0 \quad \forall x \in E$$

This is good as it means that the MCMC will be able to visit the whole support of the post distr. This property is guaranteed:
- for GS wf $\exists m \geq 1$ st $P^m$ has a (strictly) positive density $f$ (wrt $\psi$)
- for MH the one of GS plus wf $P$ has discrete and absolutely continuous components

(3) **Recurrence**. An irreducible (assumption!) MC is recurrent wf we will visit undefinitely often (wo) the sets of reachable states a unbounded amount more undefinite. Thus, or, at least from almost every virtual point.

$$P_x(X_m \in A \text{ wo}) > 0 \quad \forall x$$
$$P_x(X_m \in A \text{ wo}) = 1 \quad \psi\text{-ae wrt } x$$

**Thm** Let $(X_n)_{n\geq 0}$ an irreducible MC, and $\pi$ a stationary distr $\Rightarrow$ for wt.
- $\Rightarrow$ : the MC is $\pi$-irreducible
- : $\pi$ is the unique stat distr ( irreducible $+\ \pi$ invariant $\Rightarrow$ recurrent )
- : the MC is also recurrent

exercise

**Thm** (SLLN) Let $(X_n)_{n\geq 0}$ an irreducible MC, with $\pi$ its (unique) invariant distr. Let $f: E \to \mathbb{R}$ st $E_\pi[|f|] < \infty$. Then

$$P_x\left( \frac{1}{m+1}\sum_{\omega=0}^{m} f(X_\omega) \xrightarrow{m\to+\infty} \int_E f(x)\,\pi(x)\,dx \right) = 1 \qquad \begin{array}{l}\pi\text{-a.e.}\\ \text{w.r.t } x\end{array}$$

So starting from the "correct" point, that estimation is good. But we need to solve this initialisation problem, so that limit may not exist for $x \in C$ st $\pi(C)=0$.

Anyway the idea is really useful as to compute $\pi(A)$ we can do
$$\pi(A)=\int_A \pi(x)\,dx = \int_E [\mathbb{1}_A(x)]\,\pi(x)\,dx \implies f(x) = \mathbb{1}_A(x)$$
$$\implies \pi(A) \underset{m \text{ large}}{\approx} \frac{1}{m+1}\sum_{\omega=0}^{m} f(X_\omega) = \frac{\#(X_\omega \in A)}{m+1}$$

To solve that issue of a.e. unit $x$ we move to
(a) **Harris recurrence**. a MC is Harris recurrent wf (wt is irreducible (unskel) and then wf)
$$\forall A \text{ st } \mathbb{P}(A)>0 \ , \quad P_x(X_m \in A \text{ i.o.}) = 1 \quad \forall x \in E$$

And to be sure that we reach the invariant distr we need
(b) **aperiodicity**. a MC is aperiodic wf all states have period 1. Then we get this important

**Thm** Let $(X_n)_{n\geq 0}$ an irreducible, aperiodic MC, with $P$ wts transition matrix and $\pi$ wts invariant distr (need $\exists$, wt the MC is also recurrent). Then
$$\| P^m(x,\cdot) - \pi(\cdot) \| \xrightarrow{m\to\infty} 0 \qquad \pi\text{-a.e. wrt } x$$

Again, too, let wt be true $\forall x \in E$ we need to ask Harris recurrent instead of just recurrent.

Now small idea of the method: the goal is to approximate $E_\pi[\ell(\theta)]$, wrt the posterior distr $\pi$, and a $\ell : \Theta \to \mathbb{R}$ cord. Then we can
- build a MC $(\theta_n)_{n\geq 0}$ of state space $\Theta$ st wt is irreducible and Harris recurrent, and has $\pi(\theta|I)$ as invariant distr
- set a sensible initialisation $\theta_0$
- simulate the MC, retaine just the values after a Burn in (BI) period
- estimate
$$E_\pi[\ell(\theta)] = \int_E \ell(\theta)\,\pi(\theta|I)\,d\theta \approx \frac{1}{T+1}\sum_{\omega=0}^{T}\ell(\theta_\omega)$$

Last problem: how to set $\pi$ to be the invariant distr?
(c) **Reversibility**. a MC $(X_n)_{n\geq 0}$ of matrix $P$ is $\pi$-reversible wf
$$\pi(x)\,P(x,y) = \pi(y)\,P(y,x) \qquad \forall x \, \forall y \in E$$

**Thm** $\pi$ reversible $\Rightarrow$ $\pi$ invariant
**Proof**: we have to show $\pi P = \pi$. In the discrete case we have
$$(\pi P)_j = \sum_\omega \pi_\omega P_{\omega j} = \sum_\omega \pi_j P_{j\omega} = \pi_j\left(\sum_\omega P_{j\omega}\right) = \pi_j \cdot 1 = \pi_j$$
In the continuous case
$$\int_E \pi(x)P(x,y)\,dx = \int_E \pi(y)P(y,x)\,dx = \pi(y)\left(\int_E P(y,x)\,dx\right) = \pi(y)$$

## METROPOLIS-HASTINGS ALGORITHM

Suppose that the target distribution $\pi$ has a density (wrt a measure $\psi$). Consider a transition probability $Q(x,\tau) = g(x,\tau)$, with $g$ being the proposal density, and st $Q(x, E^+) = 1 \; \forall x \notin E^+$ where $E^+ = \mathrm{spt}(\pi) = \{x \in E : \pi(x) > 0\}$.

Then the MH alg does this:

(1) set $X_m = x$

(2) generate a candidate point $y$ sampled from $Q(x, \cdot)$      *as we need to be able to sample from this*

(3) define

$$\alpha(x,\tau) = \begin{cases} \min\left(\frac{\pi(y)\, g(y,x)}{\pi(x)\, g(x,\tau)}, 1\right) & \text{if } \mathrm{den} \neq 0 \\ 1 & \text{if } \mathrm{den} = 0 \end{cases}$$

(4) accept with probability $\alpha(x,y)$ (we generate $u \sim U([0,1])$ and decide) or reject. So we set $X_{m+1} = y$ or $X_{m+1} = x$

(5) advance to the next iteration $m+1$ and repeat

So we built a MC which has as transition kernel

$$P(x,y) = \left[\, g(x,\tau)\,\alpha(x,\tau)\,\right]\, \mathbb{1}_{\{x \neq y\}} + \left[\, r(x)\,\right]\, \mathbb{1}_{\{x = y\}} \quad \text{we } \delta_x(y)$$

$$r(x) = \mathbb{P}(x \text{ chosen in } m \to m+1) = 1 - \int_E g(x,y)\,\alpha(x,y)\,dy$$

We have to check the requirements about the MC we got, which are irreducibility, $\psi$-recurrent and $\pi$ invariant (reversible)

(1) We check that $\pi(x)$ is reversible, iff

$$\pi(x)\, \varphi(x,\tau) = \pi(y)\, \varphi(y,x)$$
$$\pi(x)\,[\, g(x,\tau)\,\alpha(x,\tau)\,] = \pi(y)\,[\, g(y,x)\,\alpha(y,x)\,]$$

Suppose we are in the case of $\alpha$ interestingly defined, ie the minimum of $\alpha$ and $1$ is $\alpha$. So $\alpha \leq 1$ and we get

$$\mathrm{LHS} = \pi(x)\,[\, g(x,y)\,\alpha(x,\tau)\,] = \pi(x)\, g(x,y)\,\frac{\pi(y)\, g(y,x)}{\pi(x)\, g(x,y)} =$$
$$= \pi(y)\, g(y,x) = \pi(y)\,[\, g(y,x)\cdot 1\,] = \mathrm{RHS}$$

   *since here we have $\alpha(y,x)$, but if $\alpha(x,y)$ was $\leq 1$ then $\alpha(y,x)$ which is $1/\alpha(x,y)$ will be $\geq 1$ so gets clipped*

(2) For irreducible and $\psi$-recurrent we need stronger assumptions on the choice of $g(x,y)$.

- Random walk MH chain: we let $g(x,y) = f(y-x)$ where $f$ is a density, like $N$. That is equivalent to setting $y = x + z$ with $z \sim f$. And we have the requirements iff $f(x) > 0 \; \forall x \in E$.

- Independence MH chain: we let $g(x,y) = f(y)$, so we are setting $y \sim f$ in an indep way from $x$. And we have the requirements iff $f(x) > 0 \; \forall$ a.e. on $E^+$.


## GIBBS SAMPLER

The MH alg becomes very inefficient when $\theta$ is multidimensional, and since it uses a joint proposal density, this can give problems when the components of $\theta$ are on different scales or multimodal or skewed.

So GS resorts to a divide & conquer approach. Let $\theta = (x,y)$ and then $\pi(x,y)$ the target distribution. Assume we know the full conditional distributions: $f_{x|y}$ and $f_{y|x}$. Then the idea is:

(1) at iteration $m$ we have $(x_m, y_m)$

(2) perform sequentially the updates:

   sample $x_{m+1}$ from $f_{x|y}(\cdot \mid y_m)$

   sample $y_{m+1}$ from $f_{y|x}(\cdot \mid x_{m+1})$

(3) advance to the next iteration and repeat

In this way we get a bivariate MC simulation, say $(X_m, Y_m)_{m \geq 0}$. Again we need to check that we have the desired/required properties on this MC. Usually these are met when:

- $\text{spt}(\pi_X) \times \text{spt}(\pi_Y) = \text{spt}(\pi)$
- $f_{X|Y}$ and $f_{Y|X}$ are $> 0$ on the respective supports of the marginals $\pi_X$ and $\pi_Y$
- the marginals $\pi_X$ and $\pi_Y$ exists (i.e. $\pi$ is not improper)

The power of this method is even clearer in higher dimensional cases, as we do

sample $X_1^{(m+1)}$ from $f_{X_1|X_2,\dots,X_p}(\cdot \mid X_2^{(m)},\dots,X_p^{(m)})$
sample $X_2^{(m+1)}$ from $f_{X_2|X_1,X_3,\dots}(\cdot \mid X_1^{(m+1)}, X_3^{(m)},\dots)$
$\vdots$
sample $X_p^{(m+1)}$ from $f_{X_p|X_1,\dots,X_{p-1}}(\cdot \mid X_1^{(m+1)},\dots,X_p^{(m+1)})$

As argued, this methods generate a MC where the transition kernel is given by the product of the full conditionals and we sample (sequentially) and always accept wt. This shows that Gibbs is a particular case of MH.

---

## SAMPLING METHODS

(1) Rejection sampling. Let $\pi(\theta)$ the target distr. from which we can't directly sample, but assume we can evaluate $\pi(\theta)$. To apply this method we need

- a proposal density $g(\theta)$, from which we are able to sample
- a constant $c$ st $\pi(\theta) \leq c \cdot g(\theta)$ $\forall \theta \in \Theta$ (i.e. an envelope)
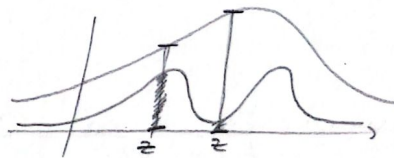
Then the method is
- draw a sample $z$ from $g(\cdot)$, and compute $r(z) = \dfrac{\pi(z)}{c \cdot g(z)}$
- generate $u \sim \mathcal{U}([0,1])$
- accept $z$ (as sample of $\pi(\cdot)$) wp $u \leq r(z)$

which is equivalent to this:
- sample $z$ from $g(\cdot)$, on the $x$ axis
- draw the vertical line till $c \cdot g(z)$
- sample $u \sim \mathcal{U}([0, c \cdot g(z)])$
- accept wp $u \leq \pi(z)$

(2) Inverse CDF. The idea: we first grab over $X$ a rv of cdf $F_X(x)$, then there is a result that $U = F_X(X) \sim \mathcal{U}([0,1])$. So

$X \sim F_X(x)$
$\Rightarrow F_X(X) \sim \mathcal{U}([0,1])$ $\Rightarrow \left( \begin{array}{c} U \sim \mathcal{U}([0,1]) \\ "U = F_X(X)" \end{array} \Rightarrow F_X^{-1}(U) \sim \mathcal{X}(x) \right)$

So the idea is, to sample from a rv $X$ of law $\mathcal{X}(x)$, of doing this:
- set $u = F_X(x)$ and invert wt, to get $x = F_X^{-1}(u)$
- sample $u_1,\dots, u_m$ from $\mathcal{U}([0,1])$
- the samples $x_i$ will be $x_i = F_X^{-1}(u_i)$

(3) Importance sampling. A method that just compute an integral, not to recall build samples from a distribution. The idea is: we are computing a $\mathbb{E}_f[\ell(\theta)]$ wrt a density $f(\cdot)$ which is difficult to sample from. So we convert to another law:

$$\mathbb{E}_f[\ell(\theta)] = \int_\Theta \ell(\theta) f(\theta)\, d\theta = \int_\Theta \ell(\theta) f(\theta) \left(\frac{g(\theta)}{g(\theta)}\right) d\theta =$$
$$= \int_\Theta \left[ \ell(\theta) \frac{f(\theta)}{g(\theta)} \right] \cdot g(\theta)\, d\theta =$$
$$= \int_\Theta \left[ \ell(\theta) w(\theta) \right] g(\theta)\, d\theta = \mathbb{E}_g[\ell(\theta) w(\theta)]$$

$g(\cdot)$ the importance distr, easier to sample from, but having the same spt of $f(\cdot)$

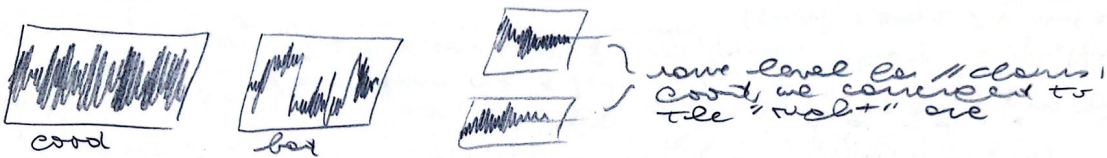(weight to correct the simulation or not among $f(\cdot)$ and $g(\cdot)$)

Small comment about the "level conditional": it was the law/
distr of a param given all the rest, we take out other params.
When computing we just leave the functional interesting part.

$$\mathcal{L}(\theta_i \mid \tau, \underline{z}) = \frac{\mathcal{L}(\theta, \tau, \underline{z})}{\mathcal{L}(\tau, \underline{z})} \propto \mathcal{L}(\theta, \tau, \underline{z}) = \mathcal{L}(\underline{z} \mid \theta, \tau) \cdot \mathcal{L}(\theta, \tau) =$$

$$= \ldots = \int_{\theta, \tau_0 \text{ wrt } \theta_i} \mathcal{L}(\underline{z} \mid \theta, \tau) \cdot \mathcal{L}(\theta) \cdot \mathcal{L}(\tau) \, \propto \ldots$$

○

---

# CONVERGENCE DIAGNOSTICS

To be checked after we do our MCMC simulation (like GS or MH)
to generate samples of the posterior distr.

(1) Trace plots. They show the history of the generated iterates for
each parameter/variable, as a time series. Is good when they look
like WN and if parallel chains settle in the end to the
same value



cood          bad

row level for // chains,
cood, we converged to
the "right" one

(2) Markov chain standard error. To gauge the variance of the estimator
$\bar{c}^{(T)}$ when we wanted to approximate the real integral $\bar{c}$.

$$\bar{c} = \int_{\Theta} c(\theta) \pi(\theta \mid \underline{z}) \, d\theta \qquad \bar{c}^{(T)} = \frac{1}{T} \sum_{\omega=1}^{T} c(\theta^{(\omega)})$$

$$\widehat{\left(\frac{\sigma^2 \bar{c}}{T}\right)} = \frac{\hat{\sigma}^2_c}{T} \left(1 + 2 \sum_{j=1}^{M} \hat{\rho}_j\right) \qquad \text{cood convergence if this value is small}$$

(3) Effective sample size. It is the # of iid iterations that we
would have to run (if it was possible) to get the same
MC standard error that we actually obtained.

So the higher is this value (the closer to our # of uncorrelated
draws) the better, as it means that our chain is "reliable".

(4) Autocorrelation plots. The bar plot of $\hat{\rho}_j = \text{Corr}(c(\theta^{(i)}), c(\theta^{(i+j)}))$.
We expect the corr to decrease as the lag increases. And the
faster it decreases the better, as it less how MC draws related
be correlated and to the previous value, all the rest are
out correlated.

# BAYESIAN LINEAR MODELS

A linear model describes a functional relation between the
mean of a response variable $\underline{y}$ and some covariates $\underline{x}$.
We have n observations, $y_1, \ldots, y_n$, and the corresponding
covariates $x_1, \ldots, x_n \in \mathbb{R}^u$. We model as

$$\boxed{y_\omega \mid x_i, \beta, \sigma^2 \overset{ind}{\sim} N(x_i^T \beta, \sigma^2) \qquad \omega = 1, \ldots, n}$$

or in matrix form:

$$\underline{y} = \begin{pmatrix} y_1 \\ \vdots \\ y_n \end{pmatrix} \quad \beta = \begin{pmatrix} \beta_0 \\ \beta_1 \\ \vdots \\ \beta_u \end{pmatrix} \in \mathbb{R}^{p=u+1}$$

$$\underline{y} \mid X, \beta, \sigma^2 \overset{\cdot}{\sim} N_n(X\beta, \sigma^2 I_n)$$

$$X = \begin{pmatrix} - x_1 - \\ \vdots \\ - x_n - \end{pmatrix} = \begin{pmatrix} 1 & x_{11} & \cdots & x_{1u} \\ \vdots & & & \\ 1 & x_{n1} & \cdots & x_{nu} \end{pmatrix}_{n \times p}$$

In general also the covariates X could be random, but if we set X and
$(\beta, \sigma^2)$ to be $\perp$ (a priori), then there will not be any functional part
involving X in the likelihood so we can discard it.

$$\mathcal{L}(\underline{y}, X \mid \omega, \beta, \sigma^2) = \mathcal{L}(\underline{y} \mid X, \omega, \beta, \sigma^2) \cdot \mathcal{L}(X \mid \omega, \beta, \sigma^2) \sim \text{const factor wrt } \beta \text{ and } \sigma^2$$

So in this context a fixed (deterministic) covariate is the mean of
a random covariate, if we assume a prior that $X \perp (\beta, \sigma^2)$

The likelihood of the model is

$$\mathcal{L}(\underline{Y}|\beta,\sigma^2) = \frac{1}{(\sigma^2)^{m/2}} e^{-\frac{1}{2\sigma^2}[(\underline{Y}-X\beta)^T(\underline{Y}-X\beta)]}$$

$$= \frac{1}{(\sigma^2)^{m/2}} e^{-\frac{1}{2\sigma^2}[s^2 + (\beta-\hat{\beta})^T X^T X (\beta-\hat{\beta})]}$$

$$\hat{\beta}_{MLE} = (X^T X)^{-1} X^T \underline{Y}$$
$$\sigma^2 = (\underline{Y}-X\hat{\beta})^T(\underline{Y}-X\hat{\beta})$$
$\Big\}$ MLE estimators for $\beta$ and $\sigma^2$ from the frequentist approach

$$B_n = (X^T X + B_0^{-1})^{-1}$$
$$b_n = B_n (X^T X \hat{\beta}_{MLE} + B_0^{-1} b_0)$$

## PRIORS AND CONJUGATE MODELS

(1) If $\sigma^2$ is known we have a conjugate prior for $\beta$.
$$\beta \sim N(b_0, B_0) \quad\Rightarrow\quad \beta|\underline{Y} \sim N(\underline{b_n}, B_n)$$
weighted average of the prior mean $b_0$ and the MLE $\hat{\beta}$

(2) Conjugate prior for $\beta$ and $\sigma^2$.
$$\pi(\beta,\sigma^2) = \pi(\beta|\sigma^2)\cdot\pi(\sigma^2)$$

$$\beta|\sigma^2 \sim N_p(b_0, \sigma^2 B_0)$$
$$\sigma^2 \sim \text{invGamma}\left(\frac{\nu_0}{2}, \frac{\nu_0\sigma_0^2}{2}\right)$$
$\Rightarrow$
$$\beta|\sigma, \underline{Y}, X \sim N_p(b_n, \sigma^2 B_n)$$
$$\sigma^2|\underline{Y}, X \sim \text{invGamma}\left(\frac{\nu_n}{2}, \frac{\nu_n\sigma_n^2}{2}\right)$$

Also here we can recover the posterior marginal of $\beta|\underline{Y}, X$, which turns out to be a multivariate (in $\mathbb{R}^p$) t student. And if we want to set predictions for new subjects, whose I know the new set of $m$ samples, we set



$$\underline{Y}_{new}|\underline{Y}, X, X_{new} \sim \\ \sim t_m(X_{new} b_n, \dots)$$

(the marginal for $\beta$ was)
$$\beta|\underline{Y}, X \sim t_p(b_n, \dots)$$

(3) Zellner's g prior. These are sets a way to set the hyperparameters of the previous model. This will is

$$\beta|\sigma^2 \sim N_p(b_0, \sigma^2 B_0)$$
$$\sigma^2 \sim \text{invGamma}\left(\frac{\nu_0}{2}, \frac{\nu_0\sigma_0^2}{2}\right)$$
$\Bigg|$
$$B_0 = c(X^T X)^{-1}$$ — this requires to be of full rank
$$\nu_0 = 0$$ — this makes the prior improper

The effect of $c$ can be seen in the posterior:
$$E[\beta|\sigma^2, \underline{Y}, X] = \frac{1}{c+1} b_0 + \frac{c}{c+1} \hat{\beta}_{MLE} \Rightarrow$$ $c$ weights the contribution of MLE vs $b_0$, the magic number is eg $\ln(m)$

(4) reference prior: the prior is improper, but the posterior comes from
$$\pi(\beta,\sigma^2) \propto \frac{1}{\sigma^2} \mathbb{1}_{(0,+\infty)}(\sigma^2) \Rightarrow$$
$$\beta|\sigma^2, \underline{Y}, X \sim N_p(\dots)$$
$$\sigma^2|\underline{Y}, X \sim \text{invGamma}(\dots)$$

(5) The most frequent / common choice is however a semi-conjugate priors, we we have a closed form for the full conditionals, and so we can use [GS].

$$\beta \sim N_p(b_0, B_0)$$
$$\tau \sim \text{Gamma}(a, b)$$
$$(\text{w/ } \beta \perp \tau)$$
$\Rightarrow$
the full conditionals are
$$\beta|\tau, \dots \sim N_p(\dots)$$
$$\tau|\beta, \dots \sim \text{Gamma}(\dots)$$

## GENERALIZED LINEAR MODELS

Here we assume that the $Y_i$ are generated from an Exp-family distr, and the mean $\mu$ of $Y$ depends on the covariates ... so $\mu = X^T\beta$ as we did before, but we now are more general, and we allow $\mu = g(X^T\beta) = g(\eta)$.

In this context we have:
- the random component, i.e. the distribution of $y_i | x_i$, which must be from the Exp. Family, and we call $\mu_i = E(y_i | x_i)$
- the linear (classical) predictor $\eta_i = x_i^T \beta$ (we'll consider later)
- the link function $g(\mu_i) = \eta_i$ (the generalization of $g(x) = x$)
- or the inverse $h(\eta_i) = \mu_i$, called response function.

The Exp-Family distr. have the form
$$ f(y_i | \theta_i, \phi) = e^{c(y_i, \phi)} \cdot e^{\left( \frac{y_i \theta_i - b(\theta_i)}{\phi} \right)} \qquad \begin{array}{l} \theta_i: \text{natural param} \\ \phi: \text{scale param} \end{array} $$

In the params we have one $\beta$ and $\phi$ (so $\theta_i$ can enter in $x_i^T \beta$). The presence of $\phi$ is according to if we model it as a random effect or not (multilevel). We usually have

cond. priors: $\pi(\beta | \phi) \cdot \pi(\phi)$ , and usually we set
$+$ priors: $\pi(\beta) \cdot \pi(\phi)$ , $\beta \sim N_p(b_0, B_0)$

Focusing on the linear regression, i.e. $y_i | x_i \stackrel{\perp}{\sim} Ber(\mu_i = g(\eta_i = x_i^T \beta))$.
We could have different choices of the link function:

$$ \mu_i = \Phi(\eta_i) \qquad\qquad \mu_i = \frac{e^{\eta_i}}{1 + e^{\eta_i}} \qquad\qquad \mu_i = 1 - e^{-e^{\eta_i}} $$

$$ [\text{probit model}] \qquad\qquad [\text{logit model}] \qquad\qquad [\text{complementary log-log model}] $$

## GIBBS SAMPLER FOR THE PROBIT MODEL
See the notes really greater, you're a rare ... can better do.
Originally, the idea was that we have the probit model
$$ y_i | x_i, \beta \stackrel{\perp}{\sim} Ber(\mu_i), \qquad \mu_i = \Phi(\eta_i) = \Phi(x_i^T \beta) $$
$$ \text{prior: } \beta \sim N_p(b_0, B_0) $$

The idea was to introduce some latent variables as
$$ z_i = x_i^T \beta + \varepsilon_i \qquad\qquad \Rightarrow \qquad y_i = \begin{cases} 1 & \text{if } z_i > 0 \\ 0 & \text{if } z_i \leq 0 \end{cases} \qquad i = 1, \dots, m $$
$$ (\varepsilon_i \stackrel{iid}{\sim} N(0,1)) $$
which is an equivalent formulation since
$$ P(y_i = 1) = P(z_i > 0) = P\left( \frac{z_i - x_i^T \beta}{1} > \frac{0 - x_i^T \beta}{1} \right) = P(Z \geq -x_i^T \beta) = $$
$$ = 1 - \Phi(-\dots) = \Phi(\dots) = \Phi(x_i^T \beta) = \mu_i $$

Then for the GS we start finding the joint law and then that the two full conditionals for $\beta$ and $z$.
$$ f(y, z, \beta) = f(y | z, \beta) \cdot f(z) \cdot f(\beta) = \underbrace{\qquad}_{} \qquad \text{as } \beta \perp z \text{ a priori} $$
$$ = \left[ \prod_{i=1}^m f(y_i | z_i) \right] \cdot \left[ \prod_{i=1}^m \underbrace{f(z_i)}_{\sim N(x_i^T \beta, 1)} \right] \cdot \underbrace{f(\beta)}_{\sim N(b_0, B_0)} $$
$$ 1_{\{y_i = 1\}} 1_{\{z_i > 0\}} + 1_{\{y_i = 0\}} 1_{\{z_i \leq 0\}} $$

and then we get the full conditionals $f(\beta | z, y)$ and $f(z | \beta, y)$.
Check the notes on this part.


# HIERARCHICAL MODELS
We now move to models which can account for latent variables, which can be complex or even impossible to measure or that were not taken into account: no variables outside the model but that should be in.

So we will have data $y$ that depends on latent variables $\gamma$ and params $\theta$, and latent will depend on params. So we have a multilevel modelling
$$ f(y | \gamma, \theta) , \quad f(\gamma | \theta) , \quad f(\theta) $$

An example of this approach (towards linear mixed effect models) is when we have clustered data, where for the levels we have the groups and the units inside the groups.

Suppose we have $J$ groups, and $Y_1, \ldots, Y_J$ with $Y_j = (Y_{1j}, \ldots, Y_{m_j j})$ representing the units inside each group. The groups are neither $\perp$ nor equally, w/ an ... so we are letting them be exchangeable. Similarly, the groups-characteristic params are not $\perp$ nor equal, so we ... then we suppose exchangeability. This leads to the multi-level model:

$$Y_{1j}, \ldots, Y_{m_j j} \mid \theta_j \overset{iid}{\sim} f(y \mid \theta_j) \qquad \text{within-group model}$$
$$\theta_1, \ldots, \theta_J \mid \phi \overset{iid}{\sim} f(\theta_j \mid \phi) \qquad \text{between-group model}$$
$$\phi \sim \pi(\phi) \qquad \text{hyper distr.}$$

Actually, we call them hierarchical or multilevel or random effects models, or the same. And the $\theta_j$ are not $\perp$ so we wish to exchange information among the groups, for example if we need to do prediction on a new group.

Taking first a above hierarchical structure, we can resort on covariates and levels

## LINEAR MIXED EFFECT MODEL

Let $Y_{wj}$ the ... corresponding to an observation, where w could be
- unit $w$ from group $j$, or
- unit $w$ measurement taken at time $j$

Then for each unit we have its vector of covariates, $x_{wj}$ for $w = 1, \ldots, m_j$ (#units in group $j$) and $j = 1, \ldots, J$ (# of groups). We want to model the connection through vectors $\beta_j$ as in a classical linear model, so we get

$$Y_{wj} = x_{wj}^T \beta_j + \varepsilon_{wj}$$
$$\varepsilon_{wj} \overset{iid}{\sim} N(0, \sigma^2)$$

or $\underset{\text{WITHIN}}{(1)}$  $\boxed{Y_j \mid X_j, \beta_j, \sigma^2 \overset{\perp}{\sim} N_{m_j}(X_j \beta_j, \sigma^2 I_{m_j})}$

This is the within-group model. For the between-group we need to decide the distr. of $\beta_j$'s. Again, they are supposed exchangeable, and we let a multilevel

$$\underset{(2)}{\text{BETWEEN}} \quad \beta_j \mid \theta, \Sigma \overset{iid}{\sim} N_p(\theta, \Sigma)$$
$$\underset{(3)}{\text{PRIOR}} \quad \theta, \Sigma \sim \pi(\theta) \cdot \pi(\Sigma)$$
$$(\text{and } \sigma^2 \sim \pi(\sigma^2))$$

This is a linear mixed effect model. The ME can be better seen through a re-parametrization.

$$\beta_j = \theta + \gamma_j \implies Y_{wj} = x_{wj}^T \beta_j + \varepsilon_{wj} =$$
$$\gamma_j \overset{iid}{\sim} N(0, \Sigma) \qquad = \boxed{x_{wj}^T \theta} + \boxed{x_{wj}^T \gamma_j} + \varepsilon_{wj}$$

$\theta$: fixed effect part, as at is constant across groups

$\gamma_j$: random effect part, as at is group specific

This can also bring the extension of different covariates close for the two parts, we:

$$Y_{wj} = x_{wj}^T \theta + z_{wj}^T \gamma_j + \varepsilon_{wj} \quad \sim \text{with } \theta \text{ and } \gamma_j \text{ also of different dim. maybe (else p and v)}$$
$$\gamma_j \mid \Sigma \overset{iid}{\sim} N_v(\theta, \Sigma)$$
$$\theta, \Sigma, \sigma^2 \sim \pi(\theta) \cdot \pi(\Sigma) \cdot \pi(\sigma^2)$$

$\llcorner$ typically we set
$$\theta \sim N_p(\psi_0, L_0)$$
$$\Sigma \sim \text{inv Wishart}(\ldots)$$
$$\sigma^2 \sim \text{inv Gamma}(\ldots)$$

# MODEL ASSESSMENT

We now use the clusters
- model selection: which model among these is the best?
- model checking: does our model fit well enough the data?

## MODEL SELECTION

(1) Compute the posterior probabilities that each model is correct
and select the model(s) with the largest probability.
- case of two models comparison: read case of the next one
- case of $k = J+1$ models comparison. We have

$$M_j: \Sigma \mid \theta_j, M_j \sim f(\Sigma \mid \theta_j, M_j), \quad \pi(\theta_j \mid M_j) = \pi(\theta_j \mid m=j) \quad j = 0 \to J$$

for $j$ in $0:J$
- let $P(m=j)$ the prior prob. of choosing model $j$;
(typically) a uniform, so $1/k$)
- compute the posterior of $\theta$ in model $M_j$:
$$\pi(\theta_j \mid \Sigma, M_j) = \frac{f(\Sigma \mid \theta_j, M_j) \cdot \pi(\theta_j \mid M_j)}{\int_{\theta_j} f(\Sigma \mid \theta_j, M_j)\,\pi(\theta_j \mid M_j)\, d\theta_j} \curvearrowright f(\Sigma \mid M_j)$$

- compute the posterior prob. mass of model $M_j$:
$$P(m=j \mid \Sigma) = \frac{f(\Sigma \mid M_j) \cdot P(m=j)}{\sum_{j=0}^{J} f(\Sigma \mid M_j) \cdot P(m=j)} \curvearrowright f(\Sigma) \leftarrow \text{like a step further down}$$
"we choose model $M_j$"

(2) Compute for each model a score about how good the model
was at predicting future observations, and choose the best-
scored one.

- KLD$_j$ (Kullback-Leibler divergence): if we knew the true law
$p(\tilde{y})$ that generates the data, we could compare wt what
the $f(\tilde{\Sigma} \mid \Sigma, M_j)$ through the KLD, for each model $M_j$.
But not knowing $p(\tilde{y})$ we have to use a proxy.

$$LPPD_j = \sum_{\omega=1}^{m} \ln\left( f(y_\omega \mid \Sigma, M_j) \right)$$

- LPPD$_j$ (log posterior predictive density): select the largest one.
But here we are using data twice, as we predict so using
all the vector $\Sigma$. So two ideas are

| removing data so when conditioning | add a new observation: |
|---|---|

$$LPML_j = \sum_{\omega=1}^{m} \ln\left( f(y_\omega \mid \Sigma_{-\omega}, M_j) \right)$$
$$= \sum_{\omega=1}^{m} \ln\left( CPO_\omega \mid M_j \right)$$
(log pseudo marginal likelihood)  conditional predictive ordinate (of unit $\omega$)

$$WAIC_j = -2(LPPD_j) + 2\,\psi w_j$$
$$\psi w_j = \sum_{\omega=1}^{m} \text{var}_{\theta_j \mid \Sigma}\left[ \ln( f(y_\omega \mid \theta_j, M_j)) \right]$$
(models applicable anywhere where)  (computable from the post values of the MCMC)

$$CPO_\omega = f(y_\omega \mid \Sigma_{(-\omega)}) = \binom{\text{posterior}}{\text{next obs}} = \int_\theta f(y_\omega \mid \theta) \cdot \pi(\theta \mid \Sigma_{(-\omega)})\, d\theta =$$

$$= \int_\theta f(y_\omega \mid \theta) \cdot \left[ \frac{\prod_{\xi=1, \xi\neq\omega}^{m} f(y_\xi \mid \theta) \cdot \pi(\theta)}{\int_\theta \prod_{\xi=1, \xi\neq\omega}^{m} f(y_\xi \mid \theta)\,\pi(\theta)\, d\theta} \right] d\theta$$

$$\Rightarrow \frac{1}{CPO_\omega} = \frac{\int_\theta \left( \prod_{\xi=1}^{m} f(y_\xi \mid \theta) \cdot \pi(\theta) \cdot \left[ f(y_\omega \mid \theta) \frac{1}{f(y_\omega \mid \theta)} \right] \right) d\theta}{\int_\theta \prod_{\xi=1}^{m} f(y_\xi \mid \theta)\,\pi(\theta)\, d\theta} =$$
$$= \pi(\theta \mid \Sigma)$$

$$= \int_\theta \frac{1}{f(y_\omega \mid \theta)}\, \pi(\theta \mid \Sigma)\, d\theta \Rightarrow \text{actually we don't need } m \text{ different MCMC samples (from each posterior } \pi(\theta_j \mid \Sigma, M_j))$$

# MODEL CHECKING

Once we have selected some models we want to be sure that they are able to correctly fit the data. We can do it with

(1) We could use the bayesian model as a data generation mechanism. From each MCMC we take e.g. the fitted values and compare their distri/quantiles to the ones of the real data.

(2) We can extend the outlier detection, to see if the real data are probable to the model we have. We call them posterior predictive probabilities:

$$p_i = \min \{ \mathbb{P}(\tilde{y} > \tilde{y}_i \mid \underline{y}, M), \; \mathbb{P}(\tilde{y} < \tilde{y}_i \mid \underline{y}, M) \} \quad i = 1, \dots, m$$

  └─ if this is too low then it is a bad news

# COVARIATE SELECTION

(1) A first idea for covariate selection is to convert this to a model selection, where we evaluate all possible models. But this scan of the entire model space becomes unfeasible when we have lots of covariate (for $n$ regressors we get $2^n$ models).

One way the idea is to select the best set of $\beta$s to build a model in the form

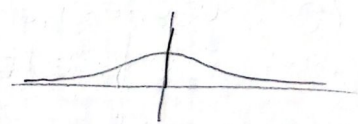$$f(y_0) = f(E[Y \mid \underline{x}, \beta]) = \eta_0 = \underline{x}^T \beta = \beta_1 x_1 + \dots + \beta_n x_n$$

Idea: we start with the full covariate models, but we assume a sparse prior on the $\beta$s to allow a regularization effect. In this way the useless covariate will get forced out.

(2) Spike and Slab. We define $\underline{\gamma} = (\gamma_1, \dots, \gamma_n)$ the vector to describe a certain model choice, as $\gamma_j \in \{0,1\}$ and $\gamma_j = 1 \Leftrightarrow \beta_j \neq 0$ (we include/remove the covariate $j$).

Then we can define a hierarchical model treating the $\beta$ as a function of $\gamma$, as $\pi(\beta, \gamma) = \pi(\beta \mid \gamma) \cdot \pi(\gamma)$. We can have
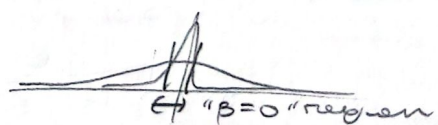
$$\beta_j \mid \gamma_j \sim (1-\gamma_j) \delta_0 + \gamma_j N(0, \sigma_j^2)$$
$$\gamma_j \mid \theta_j \sim Ber(\theta_j)$$
$$\theta_j \sim \pi(\theta_j) \quad$$
$$[\text{spike \& slab}]$$

eg $U([0,1])$, or $\theta_j = 1/2$ if to get no prior preferences



$$\beta_j \mid \sigma_j^2 \sim N(0, \sigma_j^2)$$
$$\sigma_j^2 \mid \dots, \gamma_j \sim (1-\gamma_j) \delta_{\sigma_j^2} + \gamma_j \delta_{c_0^2 \sigma_j^2}$$
$$\gamma_j \mid \theta_j \sim Ber(\theta_j)$$

[SSVS: stochastic search variable selection]



$\Leftrightarrow$ "$\beta=0$" region

And we can select the best choice by looking at the posteriors:

$$\pi(\gamma_0 \mid \underline{y}) = \frac{f(\underline{y} \mid \gamma_0) \cdot \pi(\gamma_0)}{f(\underline{y})} = \frac{f(\underline{y} \mid \gamma_0) \cdot \pi(\gamma_0)}{\sum_{\tilde{\gamma}_0} f(\underline{y} \mid \tilde{\gamma}_0) \cdot \pi(\tilde{\gamma}_0)}$$

Or if we have samples (like from a MCMC chain) we can select:

- HPD (highest posterior model/density): choose the set of $\gamma$ which occurred the most in the simulation

$$\arg\max_{\gamma_0} \frac{1}{m} \sum_{t=1}^{m} \mathbb{1}_{\{\gamma^{(t)} = \gamma_0\}}$$

- MPM (median probability model): pick all the covariates to which the posterior inclusion probability is high

$$\text{all } j: \; \pi(\gamma_j = 1 \mid \underline{y}) \approx$$
$$\approx \frac{1}{m} \sum_t \mathbb{1}_{(\gamma_j^{(t)} = 1)} > \frac{1}{2}$$

- HS (hard shrinkage): pick all the covariates to which 0 ∉ to the 95% posterior CI of $\beta_j$

$$\text{all } j: \; 0 \notin CI_{\beta_j}^{0.95}$$

# SURVIVAL ANALYSIS

Here we have data coming from studying the time until the occurrence of a certain event (an item dies, a patient dies). So the target is studying a random time span $T$. We assume
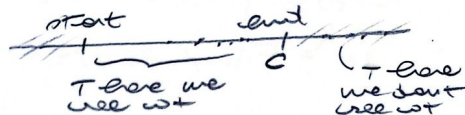
- $T \geq 0$ a.s.
- $T$ (r.v.) is absolutely continuous
- data are censored



About censoring we have

- **right censoring:** the occurring event may happen after a certain time limit, $c_i$, so we don't observe it but we just know that $T > c$. We set

$$ y_\omega = T_\omega \wedge c_\omega \qquad \delta_\omega = \begin{cases} 1 & T_\omega \leq c_\omega \\ 0 & T_i > c_\omega \end{cases} \qquad \not{y_\omega} $$

- **left censoring:** when we may have events which occurred before the start of the analysis. So we just know "they occurred before" (so the effect of right censoring) so we may set $T < t_*$.

- **interval censoring:** here we observe the data at intervals so we may get $T \in (t_*, t^*)$

Anyway, we focus on right censoring data. And we assume to have
(1) independent censoring, ie $T \perp c$ (both are r.v.)
(2) non-inormative censoring, ie the censoring distribution, wi $f_c(\cdot)$, does not depend on params of the law of $T$, $f_T(\cdot)$

To build models we need to define two functions:
- **survival function:** $S(t) = 1 - F(t) = \mathbb{P}(T > t)$
- **hazard function (or failure rate):**

$$ \ell(t) = \lim_{\delta t \to 0} \mathbb{P}(t \leq T \leq t + \delta t \mid T > t) = \frac{f(t)}{S(t)} $$

which characterizes a distribution, hence $F(t) = 1 - e^{-\int_0^t \ell(u)\, du}$

Now we suppose to have just data $(y_\omega, \delta_\omega)$ for $\omega = 1, \ldots, n$ (no covariates so now), and so we can compute the **likelihood**:

$$ (T_\omega, c_\omega) \mid f(\cdot), g(\cdot) \overset{ind}{\sim} f(\cdot) g(\cdot) $$

$$ \delta_\omega = \mathbb{1}_{\{T_\omega \leq c_\omega\}} \quad \Rightarrow \quad \varphi((\underline{y}, \underline{\delta}) \mid f, g) = \prod_{\omega=1}^{m} \varphi((y_\omega, \delta_\omega) \mid f, g) $$
$$ y_\omega = T_\omega \wedge c_\omega = T_i^{\delta_\omega} c_\omega^{1-\delta_\omega} $$

- if $\delta_\omega = 1$ (so $y_\omega = T_\omega$) $\Rightarrow \varphi(y_\omega, \delta_\omega = 1 \mid \underline\theta) = \varphi(T_\omega, T_\omega \leq c_\omega \mid \underline\theta) = $
$$ = \varphi(y_\omega, c_\omega \geq y_\omega \mid \underline\theta) = f(y_\omega)(1 - G(y_\omega)) $$

- if $\delta_\omega = 0$ (so $y_\omega = c_\omega$) $\Rightarrow \varphi(y_\omega, \delta_\omega = 0 \mid \underline\theta) = \varphi(c_\omega, T_\omega > c_\omega \mid \underline\theta) = $
$$ = \varphi(y_\omega, T_\omega > y_\omega \mid \underline\theta) = g(y_\omega)(1 - F(y_\omega)) $$

$$ \Rightarrow \varphi(\underline{y}, \underline{\delta} \mid \underline\theta) \propto \prod_{\omega=1}^{m} (f(y_\omega))^{\delta_\omega}(1 - F(y_\omega))^{1-\delta_\omega} = \qquad \text{(by non informative censoring)} $$
$$ = \prod_{\omega=1}^{m}(f(y_\omega))^{\delta_\omega}(S(y_\omega))^{1-\delta_\omega} = $$
$$ = \prod (\ell(y_\omega))^{\delta_\omega} S(y_\omega) $$

# PARAMETRIC MODELS

We want to model inference on the law of $T$, so $f(\cdot)$ & we set it be a parametric distribution, like $f_{T \mid \theta}(\cdot \mid \theta)$, and make inference on $\theta$ (like $\mathbb{E}(T \mid \theta)$, med$(T \mid \theta)$, ecc).

(1) **Exponential model:** $T_1, \ldots, T_m \mid \theta \overset{iid}{\sim} Exp(\theta)$.

$$ f(t) = \theta e^{-\theta t} \mathbb{1}_{\{t \geq 0\}} \qquad \Rightarrow \quad \ell(t) = \theta \qquad\qquad (n_m = \#\{\delta_\omega = 1\}) $$
$$ S(t) = 1 - F(t) = e^{-\theta t} $$

$$ \Rightarrow \varphi(\underline{y}, \underline{\delta} \mid \theta) \propto \prod_{\omega=1}^{m}(\ell(y_\omega))^{\delta_\omega} S(y_\omega) = (\ell(y_\omega))^{n_m} \prod S(y_\omega) $$
$$ = \theta^{n_m} \prod(e^{-\theta y_\omega}) = \theta^{n_m} e^{-\theta \sum y_i} $$

# REGRESSION MODELS (IE WITH COVARIATES)

Now as data we get $(z_w, \delta_w, x_w)$ two, still under weak censoring and the two assumptions. Being the two $z$ we can model more easily thus:

$$\ln(T_w) = x_w^T \beta + \sigma \varepsilon_w \quad , \quad \varepsilon_w \overset{iid}{\sim} F_\varepsilon \ (\text{a known distribution})$$

and we write $T_w \overset{iid}{\sim} \text{AFT}(F_\varepsilon, \beta, \sigma | x_w)$, we accelerated failure time. The name is because

$$F_T(t) = \mathbb{P}(T \geq t) = \mathbb{P}\left( \underbrace{e^{x_w^T \beta} e^{\sigma \varepsilon_w}}_{W_\sigma} \geq t \right) = \mathbb{P}(W_\sigma \geq t e^{-x_w^T \beta}) =$$

$$= F_{W\sigma}(t \cdot e^{-x_w^T \beta}) \sim \text{a cdf evaluated in a scaled time}$$

Common choices for $F_\varepsilon$

(%) $\boxed{F_\varepsilon = N(0,1)}$ in this case we get $\ln(T_w) \overset{iid}{\sim} N(x_w^T \beta, \sigma^2)$ which is called log-normal AFT. As an example this leads to:

$t^\Phi: F_T(t^\Phi) = 1/2$? (the median survival time)

$$F_T(t^\Phi) = \mathbb{P}(T \geq t^\Phi) = \mathbb{P}(\ln T \geq \ln t^\Phi) =$$
$$= \mathbb{P}\left( Z \geq \frac{\ln t^\Phi - x^T \beta}{\sigma} \right) = \bar\Phi\left( \frac{\ln t^\Phi - x^T \beta}{\sigma} \right) \overset{!}{=} \frac{1}{2}$$

$$\Rightarrow \frac{\ln t^\Phi - x^T \beta}{\sigma} = 0 \quad \Rightarrow \quad t^\Phi = e^{x^T \beta}$$

which allows to study the relative median: the ratio of $t^\Phi_1$ and $t^\Phi_2$ of two patients of equal covariates but one

$$RM(1,2) = \frac{t^\Phi_1}{t^\Phi_2} = \frac{e^{x_{11}\beta_1} \cdots e^{x_{1u}\beta_u}}{e^{x_{21}\beta_1} \cdots e^{x_{2u}\beta_u}} = e^{(x_{1u} - x_{2u})\beta_u}$$

(the all else equal covariate, ...more)

to study the impact / effect of a certain covariate.

---

# SPATIAL MODELS

In lots of contexts (environment, ecology, climate, etc) we have to work on ... which are multivariate (response and covariates), temporal and spatial.

For spatial data we divide into:
- point-referenced data (geostatistical data): where $Y(s)$ the r.v. is at location $s \in D \subset \mathbb{R}^h$ and $s$ varies continuously over $D$
- areal data: now $s \in D$ where $D$ partitioned into a entire collection of areal units, with well-defined boundaries
- point-pattern data: when $D$ is random

General idea: there results be the iid data a spatial pattern, where units closer in space and tend to be similar.

## POINT-REFERENCED DATA

We have an underlying stoch process $\{Y(s) : s \in D \subset \mathbb{R}^h = \mathbb{R}^2\}$ at $n$ locations $s_1, \ldots, s_n$. We define

- a weak stationary process
  $$\mathbb{E}[Y(s)] = y \quad \forall s$$
  $$\text{Cov}[Y(s), Y(s + \varepsilon)] = C(\varepsilon) \quad \forall s \ \forall \varepsilon$$

- a strictly stationary process
  $$\mathcal{L}(Y(s_1), \ldots, Y(s_n)) = \mathcal{L}(Y(s_1 + \varepsilon), \ldots, Y(s_n + \varepsilon))$$

- semivariogram and covariance ($C = \text{cov}$)
  $$2\gamma(\varepsilon) = \text{var}[Y(s+\varepsilon) - Y(s)]$$
  $$\gamma(\varepsilon) = C(0) - C(\varepsilon)$$

- isotropic process
  $$\gamma(\varepsilon) = \gamma(\|\varepsilon\|)$$

# GAUSSIAN REGRESSION MODEL

Regression as we have covariates at each location. We model using the following

$$Y(z) = x_{(z)}^T \beta + W(z) + \varepsilon(z)$$

a spatial residual:
(Gaussian process)
(like a random effect)

(pure (non spatial) residual:
like the model "bit" of $Y(\cdot)$

$$\varepsilon(z) \overset{iid}{\sim} N(0, \tau^2)$$
(nugget)

$$\{W(z)\} \sim GP(0, C(e) = \sigma^2 \rho(e, \varphi))$$

possible covariance models:

$$C(e) = \begin{cases} \sigma^2 e^{-\varphi e} & e > 0 \\ \sigma^2 + \tau^2 & e = 0 \end{cases}$$

$$C(e) = \begin{cases} \sigma^2 e^{-|\varphi e|^p} \\ \sigma^2 + \tau^2 \end{cases} \quad [\text{powered} \atop \text{exp}]$$

[exponential]

$$C(e) = \{ \dots \quad [\text{Matern}]$$

For for a set of locations we set $Y(z_i)$, the $x_{(z_i)}$ which will form a matrix $X$ of $m \times p$ size, and the full model is:

$$Y \mid w, \beta, \tau^2 \sim N_m(\boxed{X\beta + w}, \tau^2 I_m)$$
$$w \mid \theta \sim N_m(0, \Sigma(\theta)), \text{ where}$$
$$w = \begin{pmatrix} w_{(z_1)} \\ w_{(z_m)} \end{pmatrix} \quad \Sigma(\theta) \text{ has entries } \sigma^2 \cdot \rho(\underbrace{\|z_i - z_j\|^2}_{e}, \theta)$$

$$\beta \sim N_{4p}(\mu_\beta, \Sigma_\beta)$$
$$\theta = (\sigma^2, \varphi, \tau^2) \sim \text{informative!}$$
$$f_w) \quad (f_\varphi \quad f_\tau$$

For predictions (we bayesian theme) we can exploit the fact that $Y(z)$ is a stable (gaussian) process, so a new location so well will are $Y(z_0)$ gaussian.

$$f(Y_0 \mid Y, X, z_0) = \int f(Y_0 \mid Y, \beta, \theta, z_0) \cdot f(\beta, \theta \mid Y, X) \, d\beta \, d\theta$$

$$\begin{pmatrix} Y_0 \\ Y \end{pmatrix} \in \mathbb{R}^{m+4} \mid \beta, \theta \sim N_{m+4}\left(\tilde{X}\beta = \left(\frac{\tilde{x}_0}{X}\right)\beta, \tilde\Sigma_{augmented}\right)$$

# AREAL DATA

Here we start defining a adjacency matrix $W = [w_{ij}]$, which tells us if locations $i$ and $j$ are close, and if $w_{ij} = 1$ if then (areas share a boundary) or not, and if $w_{ij} = 0$.

This will allow us to treat/describe more easily the full conditional laws, as we will then use $f(Y_0 \mid Y_j : j \in \partial w)$, with $\partial w$ the set of neighbours of $w$.

From read/using these full conditionals actually we are able to recover the joint law of $(Y_1 \dots Y_m)$. This procedure is called MRF (Markov Random Field).

The CAR model (conditional auto-regressive) is an example of MRF for the conditional joint is a full distr (it not exists but can be computed). We model

$$\text{(the full conditionals)} \quad Y_i \mid Y_{(-i)} \sim N\left(\frac{\sum_j w_{ij} Y_j}{\sum_j w_{ij}}, \frac{\tau_i^2}{\sum_j w_{ij}}\right)$$

$$\Rightarrow \text{(the joint distr)} \quad p(Y_1 \dots Y_m) \propto e^{[-\frac{1}{2\tau^2} Y^T(D_W - W)Y]}$$

$$D = \begin{pmatrix} \tau_1^2 & & 0 \\ & \ddots & \\ 0 & & \tau_m^2 \end{pmatrix} \quad W = \begin{pmatrix} w_1. & & 0 \\ & \ddots & \\ 0 & & w_m. \end{pmatrix}$$

(this matrix actually is singular, so the when is singular (as we cant invert it)

solution: we $\Sigma_i^{-1} = D_W - \rho W$
- choose a $\rho$ to make it non singular
- put a prior on $\rho$ ~ $N(0,1)$