

(Generalized) Linear Mixed Effect models

Alessandra Guglielmi

Politecnico di Milano
Dipartimento di Matematica
Milano, Italia
e-mail: alessandra.guglielmi@polimi.it

November 18, 2020



A. Guglielmi

Bayesian Statistics

1

Bayesian Hierarchical Models (from the first lesson)

Multilevel data: results of a test for students in a population of schools in US

- test/exam scores of students in different schools or universities
- failure times of items in different batches
- patients within several hospitals
- people (or items) within provinces within regions within countries

Two levels:

- groups
- units within groups

y_{ij} is the data of the i -th unit in group j : $i = 1, \dots, n_j$,
 $j = 1, \dots, m$

$$(\mathbf{Y}_1, \mathbf{Y}_2, \dots, \mathbf{Y}_m) \quad \mathbf{Y}_j = (Y_{1,j}, \dots, Y_{n_j,j})$$

A. Guglielmi

Bayesian Statistics

2

Bayesian Gaussian hierarchical model (ANOVA)

We considered:

$$\begin{aligned} Y_{1,j}, \dots, Y_{n_j,j} | \theta_j &\stackrel{iid}{\sim} N(\theta_j, \sigma^2) && \text{group-specific parameters} \\ \theta_1, \dots, \theta_m | (\mu, \tau^2) &\stackrel{iid}{\sim} N(\mu, \tau^2) && \text{within-group model} \\ &(\mu, \tau^2) \sim \pi && \text{between-group model} \end{aligned}$$

we assume them to be exchangeable (conditioned on (μ, τ^2))

we assume the variance to be equal in all groups: what vary it's the mean score for each group

population of groups/group-parameters: prediction on a student coming from a new school, selected at random from the population of groups

- group-specific parameters (each θ_j is the mean of the math score in the j -th school)
- $\theta_1, \dots, \theta_m$ are NOT independent, since we want to share information between the groups; the dependency is *mild* (exchangeability)

A. Guglielmi

Bayesian Statistics

3

Exchangeable prior for $\theta_1, \dots, \theta_m$

exchangeability of the prior of $\theta_1, \dots, \theta_m$ allows the groups to exchange information (between the groups)

- We do not possess any knowledge about the groups
- all groups are considered similarly, but NOT independent
- thanks to exchangeability of the prior for $(\theta_1, \dots, \theta_m)$, groups with few data borrow strength from larger groups
- Bayesian inference for any θ_j reflects a combination of info about θ_j from the data in group j and the hierarchical component of the model $\theta_j \sim N(\mu, \tau^2)$
- data are used to update prior beliefs about (μ, τ^2)

As a result, data from group j helps shape the posterior over θ_k ($\forall k \neq j$) via its contribution to inference for (μ, τ^2)

sharing/borrowing information across groups



A. Guglielmi

Bayesian Statistics

4

Borrowing strength (first lesson)

$$E(\theta_j | \bar{y}_j, \mu, \tau^2, \sigma^2) = \frac{n_j/\sigma^2}{n_j/\sigma^2 + 1/\tau^2} \bar{y}_j + \frac{1/\tau^2}{n_j/\sigma^2 + 1/\tau^2} \mu$$

frequentist estimator of θ_j prior mean of θ_j

When n_j is small, i.e. group j gives little info about θ_j , i.e. the frequentist estimate \bar{y}_j is poor; however the Bayesian estimate:

$$E(\theta_j | \dots) \approx \mu$$

The Bayesian estimate is obtained borrowing strength from the other groups (through μ)



Borrowing strength (first lesson)

$$E(\theta_j | \bar{y}_j, \mu, \tau^2, \sigma^2) = \frac{n_j/\sigma^2}{n_j/\sigma^2 + 1/\tau^2} \bar{y}_j + \frac{1/\tau^2}{n_j/\sigma^2 + 1/\tau^2} \mu$$

frequentist estimator of θ_j prior mean of θ_j

When τ^2 is large (heterogeneous groups), the Bayesian estimate:

$$E(\theta_j | \dots) \approx \bar{y}_j$$

there is less shrinkage to μ , relying more on the info in group j .



Bayesian Gaussian hierarchical model (first lesson)

The prior is completed assuming:

$$\begin{aligned} \frac{1}{\sigma^2} &\sim \text{gamma}\left(\frac{\nu_0}{2}, \frac{\nu_0 \sigma_0^2}{2}\right) & \sigma^2 &\text{within-group variance} \\ \frac{1}{\tau^2} &\sim \text{gamma}\left(\frac{\eta_0}{2}, \frac{\eta_0 \tilde{\sigma}_0^2}{2}\right) & \tau^2 &\text{between-group variance} \\ \mu &\sim N(\mu_0, \gamma_0^2) \end{aligned}$$

R Example: Bayesian_hierarchical_primalezione.R



Math score example + covariate: least squares estimates within groups

Suppose now that for every student we have an extra information: SES_{ij}

Math score example, $m = 100$ schools
 Y_{ij} math score, SES_{ij} socioeconomic status of student i in school j
(economy + education of the family)

Least squares estimates of the m different models, for $j = 1, \dots, m$:

$$\mathbb{E}(Y_{ij}) = \beta_{0j} + \beta_{1j} \times SES_{ij}, \quad i = 1, \dots, n_j$$

$$(\hat{\beta}_{0j}, \hat{\beta}_{1j}), j = 1, \dots, m$$

What about a Bayesian model to analyze these data?



Linear models with mixed effects

y_{ij} = response of unit i in group j : $i = 1, \dots, n_j, j = 1, \dots, m$

($\mathbf{Y}_1, \mathbf{Y}_2, \dots, \mathbf{Y}_m$), where each $\mathbf{Y}_j := (Y_{1,j}, \dots, Y_{n_j,j})$

(\mathbf{x}_{ij}) a p -dim vector of covariates of subject i in group j (in this case it's the SES $_{ij}$)

(\mathbf{X}_j) the $n_j \times p$ matrix of covariates in group j

(β_j) a p -dim vector of regression parameters, for each j (\leftarrow it's a vector : $\vec{\beta}_j$)

Likelihood:

$$Y_{ij} = \mathbf{x}_{ij}^T \beta_j + \epsilon_{ij}, \quad \epsilon_{ij} \stackrel{\text{i.i.d.}}{\sim} \mathcal{N}(0, \sigma^2) \quad (1)$$

Equivalently:

$$Y_j \stackrel{\text{ind.}}{\sim} \mathcal{N}_{n_j}(\mathbf{X}_j \beta_j, \sigma^2 I_{n_j})$$

β_1, \dots, β_m are the regression parameters within the groups:
which prior? (↓)



Exchangeable prior for the within-group regression parameters

- if a priori β_1, \dots, β_m are independent $\Rightarrow \beta_1, \dots, \beta_m$ are independent a posteriori as well, i.e. same inference as m distinct Bayesian models (one within each group)
- $\beta_1 = \dots = \beta_m \Rightarrow$ analysis within one single group

Intermediate prior: $\pi(\beta_1, \dots, \beta_m)$ exchangeable, i.e.

$$\beta_j | \theta, \Sigma \stackrel{\text{i.i.d.}}{\sim} \mathcal{N}_p(\theta, \Sigma)$$

$$\theta, \Sigma \sim \pi(\theta, \Sigma)$$

hierarchical prior (weak information)

very easy to compute the predictive distribution of a new student coming from a new school

$$(\beta_j = \vec{\beta}_j)$$

we assume a prior for $\vec{\theta}$ and Σ so that the marginal joint distribution for all the betas will be exchangeable.

Why exchangeability?

First of all because it's simple, but more than that because in this way groups can exchange informations. Moreover, in this case it's very easy to compute the predictive distribution of a new student (in general new element), even coming from a new school (in general new group)!

Linear mixed effects model - LMM

LMM model

Summing up, for $i = 1, \dots, n_j, j = 1, \dots, m$:

$$Y_{ij} = \mathbf{x}_{ij}^T \beta_j + \epsilon_{ij}, \quad \epsilon_{ij} \stackrel{\text{i.i.d.}}{\sim} \mathcal{N}(0, \sigma^2)$$

$$\beta_j | \theta, \Sigma \stackrel{\text{i.i.d.}}{\sim} \mathcal{N}_p(\theta, \Sigma) \quad j = 1, \dots, m$$

$$\theta, \Sigma \sim \pi(\theta, \Sigma) = \pi(\theta) \times \pi(\Sigma)$$

$$\sigma^2 \sim \pi(\sigma^2)$$

a prior we assume $\vec{\theta} \perp\!\!\!\perp \Sigma$



Alternative parameterization

to explain
the name:
"mixed effect
model"

If $\beta_j = \theta + \gamma_j$, then (1) gives:

$$Y_{ij} = \mathbf{x}_{ij}^T \beta_j + \epsilon_{ij} = \mathbf{x}_{ij}^T \theta + \mathbf{x}_{ij}^T \gamma_j + \epsilon_{ij}, \quad \epsilon_{ij} \stackrel{\text{i.i.d.}}{\sim} \mathcal{N}(0, \sigma^2)$$

and the hierarchical prior is:

$$\gamma_j | \Sigma \stackrel{\text{i.i.d.}}{\sim} \mathcal{N}_p(\mathbf{0}, \Sigma)$$

θ : population mean effect

θ : fixed effect parameter, γ_j : random effect parameter (varies between groups)

The model contains both fixed effect and random effect parameters, so that it is called mixed effects model

Remark: here $\dim(\theta) = \dim(\gamma_j)$.



Linear mixed effect (hierarchical) models

More generally:

$$Y_{ij} = \mathbf{x}_{ij}^T \boldsymbol{\theta} + \mathbf{z}_{ij}^T \boldsymbol{\gamma}_j + \epsilon_{ij}, \quad \epsilon_{ij} \stackrel{\text{i.i.d.}}{\sim} \mathcal{N}(0, \sigma^2)$$

$$\dim(\mathbf{x}_{ij}) = \dim(\boldsymbol{\theta}) = p \quad \dim(\mathbf{z}_{ij}) = \dim(\boldsymbol{\gamma}_j) = v$$

$$\mathbf{z}_{ij} | \Sigma \stackrel{\text{i.i.d.}}{\sim} \mathcal{N}_v(\mathbf{0}, \Sigma) \quad j = 1, \dots, m$$

$$\boldsymbol{\theta}, \Sigma, \sigma^2 \sim \pi(\boldsymbol{\theta}, \Sigma, \sigma^2) = \pi(\boldsymbol{\theta}) \times \pi(\Sigma) \times \pi(\sigma^2)$$

it's not necessary that
 $\dim(\boldsymbol{\theta}) = \dim(\boldsymbol{\gamma}_j)$

Typically:

$$\boldsymbol{\theta} \sim \mathcal{N}_p(\boldsymbol{\mu}_0, L_0)$$

$$\Sigma \sim \text{inverse-Wishart}(S_0^{-1}, \nu_0)$$

$$\sigma^2 \sim \text{inverse-gamma}\left(\frac{\nu_0}{2}, \frac{\nu_0 \sigma_0^2}{2}\right)$$

Linear mixed effect (hierarchical) models

- Often $\mathbf{z}_{ij} = \mathbf{x}_{ij}$, so that $\boldsymbol{\gamma}_j$ is parameter representing group-specific effect on the covariates \mathbf{x}_{ij}
- Often $\mathbf{z}_{ij} = \mathbf{z}_j$, e.g. dummy variables representing the subgroup, or group-specific covariates

Ex 7.10 in Jackman's textbook:

y_{ij} : maths score of student i in school j

x_{ij} : ses of student i in school j

$z_{j1}=1$ if school j is a Catholic school, =0 otherwise → only present if the school is catholic

z_{j2} : average ses level in school j

$$E(y_{ij} | x_{ij}, (z_{j1}, z_{j2})) = \gamma_{10} + \gamma_{11}z_{j1} + \gamma_{12}z_{j2} + (\gamma_{20} + \gamma_{21}z_{j1} + \gamma_{22}z_{j2})x_{ij}$$

$$= \beta_{j1} + \beta_{j2}x_{ij}$$

effect of the average SES

PRIOR:

$$(\beta_{j1}, \beta_{j2})^t \sim \mathcal{N}_2((\gamma_{10} + \gamma_{11}z_{j1} + \gamma_{12}z_{j2}, \gamma_{20} + \gamma_{21}z_{j1} + \gamma_{22}z_{j2})^t, \Omega)$$

Some remarks

- ① the components of \mathbf{Y}_j , conditionally on $(\boldsymbol{\theta}, \boldsymbol{\gamma}_j)$, are independent, i.e. $\text{Var}(\mathbf{Y}_j | \boldsymbol{\theta}, \boldsymbol{\gamma}_j) = \sigma^2 I_{n_j}$, but marginally w.r.t. $\boldsymbol{\gamma}_j$, they are NOT independent:

when we marginalize we lose the independence

$$\mathbf{Y}_j | \boldsymbol{\theta}, \Sigma, \sigma^2 \sim \mathcal{N}_{n_j}(\mathbf{X}_j \boldsymbol{\theta}, \mathbf{Z}_j \Sigma \mathbf{Z}_j^T + \sigma^2 I_{n_j})$$

- ② Generalized linear mixed models

we considered the likelihood to be gaussian but we can generalize to generalized linear mixed models

Predictive distributions for a new subject from a new group

We can simulate from the predictive distribution:

- ③ Prediction for a new subject from a new group:

$$\begin{aligned} & \mathcal{L}(Y_{m+1}^{\text{new}}, \beta_{m+1} | \mathbf{y}_1, \dots, \mathbf{y}_m, \mathbf{X}_1, \dots, \mathbf{X}_m, \mathbf{x}^{\text{new}}) \\ &= \int \mathcal{L}(Y_{m+1}^{\text{new}}, \beta_{m+1}, d\beta_1, \dots, d\beta_m, d\theta, d\Sigma | \mathbf{y}_1, \dots, \mathbf{y}_m, \mathbf{X}_1, \dots, \mathbf{X}_m, \mathbf{x}^{\text{new}}) \\ &= \int \mathcal{L}(Y_{m+1}^{\text{new}} | \beta_{m+1}, \mathbf{x}^{\text{new}}) \mathcal{L}(d\beta_1, \dots, d\beta_m, \beta_{m+1}, d\theta, d\Sigma | \text{data}, \mathbf{x}^{\text{new}}) \\ &= \int \mathcal{L}(Y_{m+1}^{\text{new}} | \beta_{m+1}, \mathbf{x}^{\text{new}}) \mathcal{L}(\beta_{m+1} | \beta_1, \dots, \beta_m, \theta, \Sigma, \text{data}, \mathbf{x}^{\text{new}}) \end{aligned}$$

$$\times \mathcal{L}(d\theta, d\Sigma, d\beta_1, \dots, d\beta_m | \text{data})$$

$$= \int \mathcal{L}(Y_{m+1}^{\text{new}} | \beta_{m+1}, \mathbf{x}^{\text{new}}) \mathcal{L}(\beta_{m+1} | \theta, \Sigma) \mathcal{L}(d\theta, d\Sigma, d\beta_1, \dots, d\beta_m | \text{data})$$

$$\text{since } \mathcal{L}(\beta_{m+1} | \beta_1, \dots, \beta_m, \theta, \Sigma, \text{data}, \mathbf{x}^{\text{new}}) =$$

$$\mathcal{L}(\mathbf{Y} | \beta_1, \dots, \beta_m, \mathbf{X}) \mathcal{L}(\beta_1, \dots, \beta_m, \beta_{m+1} | \theta, \Sigma) \mathcal{L}(\theta, \Sigma) \propto$$

$$\mathcal{L}(\beta_{m+1} | \theta, \Sigma)$$

posterior from which we should be able to sample from group specific distribution $(\mathcal{N}(\boldsymbol{\theta}, \Sigma))$ (slide 10)

What we need to do is sample from the joint predictive distribution of \mathbf{Y}^{new} joint with $\boldsymbol{\beta}^{\text{new}}$ through MCMC. In the end we'll discard $\boldsymbol{\beta}^{\text{new}}$ and we'll consider only values of \mathbf{Y}^{new}

likelihood of the data

❸ Prediction for a new subject from an existing group:

$$\begin{aligned} \mathcal{L}(Y_j^{\text{new}} | \mathbf{y}_1, \dots, \mathbf{y}_m, X_1, \dots, X_m, \mathbf{x}^{\text{new}}) \\ = \int \mathcal{L}(Y_j^{\text{new}}, d\beta_j, d\theta, d\Sigma | \mathbf{y}_1, \dots, \mathbf{y}_m, X_1, \dots, X_m, \mathbf{x}^{\text{new}}) \\ = \int \mathcal{L}(Y_j^{\text{new}} | \beta_j, \mathbf{x}^{\text{new}}) \mathcal{L}(d\beta_j, d\theta, d\Sigma | \text{data}) \end{aligned}$$



Gelman (2006)- Noninformative prior for σ_0

ANOVA Model:

overall population effect

$$y_{ij} | \text{par} \stackrel{\text{ind.}}{\sim} \mathcal{N}(\mu + \alpha_j, \sigma_y^2), i = 1, \dots, n_j, j = 1, \dots, m$$

Options: $\alpha_j | \sigma_\alpha \stackrel{\text{i.i.d.}}{\sim} \mathcal{N}(0, \sigma_\alpha^2), j = 1, \dots, m$
 $\sigma^2 \sim \text{inv-gamma}(\epsilon, \epsilon)$ or $\sigma_\alpha \sim \mathcal{U}(0, \sigma_0)$

group specific effect

e.g. $\epsilon = 0.01, \sigma_0 = 100$

Gelman (2006) argues that:

- $\sigma^2 \sim \text{inv-gamma}(\epsilon, \epsilon)$ makes posterior inference sensitive to ϵ
- better $\sigma_\alpha \sim \mathcal{U}(0, \sigma_0)$ or $\sigma_\alpha \sim \text{half-t family}$

Remark: A priori $E(\alpha_j) = 0$ for all j ; otherwise the model is not identifiable



Covariance matrix distributions

A matrix M is positive definite if $\mathbf{z}^T M \mathbf{z} > 0$ for all $\mathbf{z} \neq \mathbf{0}$ (\Leftrightarrow tutti i minori principali della matrice sono > 0)

Definition: A $p \times p$ matrix X , symmetric and positive definite, has Wishart density with parameter (M, ν) , where $\nu \geq p$ and M symmetric and positive definite $p \times p$ matrix :

$$f(X | M, \nu) = \frac{1}{2^{\frac{\nu p}{2}} \Gamma_p(\frac{\nu}{2}) |M|^{\nu/2}} |X|^{\frac{\nu-p-1}{2}} \exp \left\{ -\text{tr} \left(\frac{M^{-1} X}{2} \right) \right\},$$

where $\Gamma_p(\frac{\nu}{2}) = \pi^{p(p-1)/4} \prod_{j=1}^p \Gamma(\frac{\nu+1-j}{2})$ and
 $\text{tr}(A) = \sum_j a_{jj}$ (sum of the diagonal elements of the matrix)

$$X \sim \text{Wishart}(M, \nu)$$



Alternative definition of Wishart distribution

$$\mathbf{Y}_i \stackrel{\text{i.i.d.}}{\sim} \mathcal{N}_p(\mathbf{0}, \Sigma), i = 1, \dots, n, \quad \Sigma \text{ positive definite}$$

Build the (positive definite) matrix

$$X = \sum_{i=1}^n \mathbf{Y}_i \mathbf{Y}_i^T \sim \text{Wishart}(\Sigma, n)$$

Ex: $p = 2$

$$\begin{bmatrix} \sum_{i=1}^n Y_{i1}^2 & \sum_{i=1}^n Y_{i1} Y_{i2} \\ \sum_{i=1}^n Y_{i1} Y_{i2} & \sum_{i=1}^n Y_{i2}^2 \end{bmatrix}$$



Covariance matrix distributions

✓ If $X \sim \text{Wishart}(M, \nu)$, then

$$\mathbb{E}(X|M, \nu) = \nu M$$

$$\text{Var}(X_{ij}|M, \nu) = \nu(m_{ij}^2 + m_{ii}m_{jj}), \quad \text{where } X = [X_{ij}]_{ij}, M = [M_{ij}]_{ij}$$

✓ It is a multivariate version of the gamma distribution; in fact, if $p = 1$ and $M = 1$, then $X \sim \chi^2(\nu) = \text{gamma}(\nu/2, 1/2)$

✓ the Wishart density is the conjugate prior for precision matrix Σ^{-1} under multivariate normal likelihood:

$$Y_i \sim N_p(\mu_0, \Sigma), \quad \Sigma^{-1} \sim \text{Wishart}(M, \nu) \\ \Rightarrow \pi(\Sigma^{-1} | \mathbf{y}) \sim \text{Wishart}(\cdot, \cdot)$$



Inverse Wishart distribution for a matrix X

Definition: A $p \times p$ matrix W , symmetric and positive definite, has inverse-Wishart(M^{-1}, ν) if $W^{-1} \sim \text{Wishart}(M, \nu)$.

✓ If W , a $p \times p$ matrix W , symmetric and positive definite,

$$W \sim \text{inverse-Wishart}(M^{-1}, \nu),$$

then

$$\mathbb{E}(W) = \frac{1}{\nu - p - 1} M^{-1}, \quad \nu > p + 1.$$

Be careful of the notation from software to software



Example in Ch 11 - P. Hoff's book

For $i = 1, \dots, n_j, j = 1, \dots, m$:

$$Y_{ij} = \mathbf{x}_{ij}^T \boldsymbol{\beta}_j + \epsilon_{ij} = \beta_{1j} + \beta_{2j} SES_{ij} + \epsilon_{ij}, \quad \epsilon_{ij} \stackrel{\text{i.i.d.}}{\sim} N(0, \sigma^2) \\ \beta_1, \dots, \beta_m | \boldsymbol{\theta}, \Sigma \stackrel{\text{i.i.d.}}{\sim} \mathcal{N}_2(\boldsymbol{\theta}, \Sigma), \quad (\begin{matrix} \beta_1 \\ \vdots \\ \beta_m \end{matrix}) \\ \boldsymbol{\theta} \sim \mathcal{N}_2(\boldsymbol{\mu}_0, L_0) \\ \Sigma \sim \text{inverse-Wishart}(S_0^{-1}, \eta_0) \\ \sigma^2 \sim \text{inverse-gamma}\left(\frac{\nu_0}{2}, \frac{\nu_0 \sigma_0^2}{2}\right)$$

Parameters: $\beta_1, \dots, \beta_m, \boldsymbol{\theta}, \Sigma, \sigma^2$

$\dim(\beta_i) = \dim(\boldsymbol{\theta}) = 2; \Sigma: 2 \times 2 \text{ matrix}; \sigma^2 > 0$



Example in Ch 11 - P. Hoff's book

In particular:

$\boldsymbol{\mu}_0$ =mean of the Least Squares estimates within groups

L_0 =sample covariance of the LSEs (weak prior information), i.e. a priori 95% CI for the slope parameter θ_2 is (-3.86, 8.60) that is quite large given that SES is centered and varies in (-2.5, 2.5)

$\nu_0 = 1, \sigma_0^2$ =average of the within-group sample variances

$\eta_0 = p + 2 = 4, S_0 = L_0$: the prior for Σ is diffuse, but

$\mathbb{E}(\Sigma)$ =sample covariance of the LS estimates

Remark: According to Hoff notation, $\mathbb{E}(\Sigma) = \frac{1}{\eta_0 - p - 1} S_0$



Models for longitudinal data

In some applications, j is the index of the individuals in a sample and $i = 1, \dots, n_j$ represents repeated measurements of the same variable for the same individual j , under different experimental conditions or at different (discrete) times t_{ij} (e.g., $t_{ij} = t_i$, all individuals are evaluated at the same time periods)

$\mathbf{y}_j = (y_{1j}, \dots, y_{n_j j})$ is the vector of all measurements, or the time series, of subject j in the sample

Assuming a LMM (or a GLMM) for $\mathbf{y}_1, \dots, \mathbf{y}_m$ implies we assume the components of each \mathbf{y}_j dependent

Typically, we model $E(Y_{ij}|par)$; in LMMs:

$$Y_{ij} | \mu_{ij} \sim \mathcal{N}(\mu_{ij}, \sigma^2)$$

$$\mu_{ij} = \mathbf{x}_{ij}^T \boldsymbol{\theta} + \mathbf{z}_{ij}^T \boldsymbol{\gamma}_j$$



Examples of GLMMs for longitudinal data

Linear structure for the mean: $\mu_{ij} = \theta_0 + \theta_1 t_i + \gamma_{0j} + \gamma_{1j} t_i$

θ_0 : population intercept

γ_{0j} : deviation of the intercept of the j -th subject to the popul. intercept

θ_1 : population slope

γ_{1j} : deviation of the slope of the j -th subject to the popul. slope

Each subject j has a different intercept (for the mean) at time t_i equal to $\theta_0 + \gamma_{0j}$ and a different slope equal to $\theta_1 + \gamma_{1j}$

An alternative is to model jointly the joint distribution of the vector of all the measurements of a subject (for instance through an autoregressive model or an ARMA model,

which means that we're assuming that the response at time i is related (linearly) to the response at time $i-1$ (plus a random error)

AR(1) models: $Y_{ij} = \alpha_j Y_{i-1j} + \varepsilon_{ij}$, $\varepsilon_{ij} \stackrel{i.i.d.}{\sim} \mathcal{N}(0, \sigma^2)$
it allows for the inclusion of time-dependent covariates \mathbf{x}_{ij}
If $Y_1 \sim \mathcal{N}(\dots, \sigma^2)$, then $Y_j \sim \mathcal{N}_{\eta_j}(\dots, V_j)$, where (here $\eta_j = 4$)

$$V_j = \begin{pmatrix} \sigma^2 & \alpha_j & \alpha_j^2 & \alpha_j^3 \\ \sigma^2 & \alpha_j & \alpha_j^2 & \alpha_j^3 \\ \sigma^2 & \alpha_j & \alpha_j^2 & \alpha_j^3 \\ \sigma^2 & \alpha_j & \alpha_j^2 & \alpha_j^3 \end{pmatrix}$$



Exchangeable prior for $(\alpha_1, \dots, \alpha_m)$

very simple!

(good advantage)

we need only 2 parameters

to represent V_j (we don't need the inverse-Wishart & co.)

11.3 Posterior analysis of the math score data

To analyze the math score data we will use a prior distribution that is similar in spirit to the unit information priors that were discussed in Chapter 9. For example, we'll take μ_0 , the prior expectation of θ_1 to be equal to the average of the ordinary least squares regression estimates and the prior variance Λ_0 to be their sample covariance. Such a prior distribution represents the information of someone with unbiased but weak prior information. For example, a 95% prior confidence interval for the slope parameter θ_2 under this prior is (-3.86, 8.60), which is quite a large range when considering what the extremes of the interval imply in terms of average change in score per unit change in SES score. Similarly, we will take the prior sum of squares matrix \mathbf{S}_0 to be equal to the covariance of the least squares estimate, but we'll take the prior degrees of freedom η_0 to be $p + 2 = 4$, so that the prior distribution of Σ is reasonably diffuse but has an expectation equal to the sample covariance of the least squares estimates. Finally, we'll take σ_0^2 to be the average of the within-group sample variance but set $\nu_0 = 1$.

example, the first plot in Figure 11.3 shows the posterior distribution of θ_2 , the expected within-school slope parameter. A 95% quantile-based posterior confidence interval for this parameter is (1.83, 2.96), which, compared to our prior interval of (-3.86, 8.60), indicates a strong alteration in our information about θ_2 .

The fact that θ_2 is extremely unlikely to be negative only indicates that the population average of school-level slopes is positive. It does not indicate that any given within-school slope cannot be negative. To clarify this distinction, the posterior predictive distribution of $\hat{\beta}_2$, the slope for a to-be-sampled school, is plotted in the same figure. Samples from this distribution can be generated by sampling a value $\hat{\beta}_2^{(s)}$ from a multivariate normal($\boldsymbol{\theta}^{(s)} ; \Sigma^{(s)}$) distribution for each scan s of the Gibbs sampler. Notice that this posterior predictive distribution is much more spread out than the posterior distribution of θ_2 , reflecting the heterogeneity in slopes across schools. Using the Monte Carlo approximation, we have $\Pr(\hat{\beta}_2 < 0 | \mathbf{y}_1, \dots, \mathbf{y}_m, \mathbf{X}_1, \dots, \mathbf{X}_m) \approx 0.07$, which is small but not negligible.

The second panel in Figure 11.3 plots posterior expectations of the 100 school-specific regression lines, with the line given by the posterior mean of $\boldsymbol{\theta}$ in black. Comparing this to the first panel of Figure 11.1 indicates how the hierarchical model is able to share information across groups, shrinking extreme regression lines towards the across-group average. In particular, hardly any of the slopes are negative when we share information across groups.

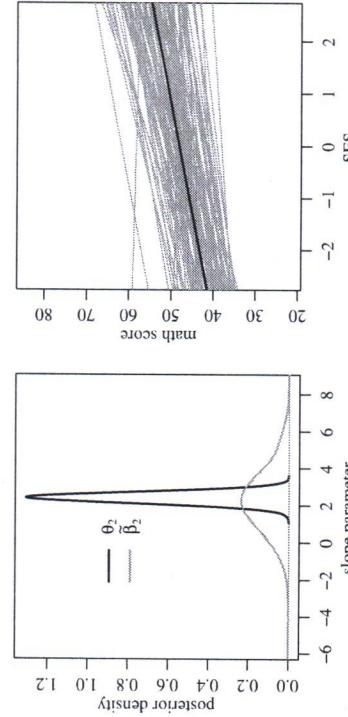


Fig. 11.3. Relationship between SES and math score. The first panel plots the posterior density of the expected slope θ_2 of a randomly sampled school, as well as the posterior predictive distribution of a randomly sampled slope. The second panel gives posterior expectations of the 100 school-specific regression lines, with the average line given in black.

Running a Gibbs sampler for 10,000 scans and saving every 10th scan produces a sequence of 1,000 values for each parameter, each sequence having a fairly low autocorrelation. For example, the lag-10 autocorrelations of θ_1 and θ_2 are -0.024 and 0.038. As usual, we can use these simulated values to make Monte Carlo approximations to various posterior quantities of interest. For

11.4 Generalized linear mixed effects models

As the name suggests, a generalized linear mixed effects model combines aspects of linear mixed effects models with those of generalized linear models described in Chapter 10. Such models are useful when we have a hierarchical data structure but the normal model for the within-group variation is not appropriate. For example, if the variable Y were binary or a count, then more appropriate models for within-group variation would be logistic or Poisson regression models, respectively.

A basic generalized linear mixed model is as follows:

$$\begin{aligned} \boldsymbol{\beta}_1, \dots, \boldsymbol{\beta}_m &\sim \text{i.i.d. multivariate normal}(\boldsymbol{\theta}, \Sigma) \\ p(\mathbf{y}_j | \mathbf{X}_j, \boldsymbol{\beta}_j, \gamma) &= \prod_{i=1}^{n_j} p(y_{i,j} | \boldsymbol{\beta}_j^T \mathbf{x}_{i,j}, \gamma), \end{aligned}$$

with observations from different groups also being conditionally independent. In this formulation $p(y | \boldsymbol{\beta}^T \mathbf{x}, \gamma)$ is a density whose mean depends on $\boldsymbol{\beta}^T \mathbf{x}$, and γ is an additional parameter often representing variance or scale. For example, in the normal model $p(y | \boldsymbol{\beta}^T \mathbf{x}, \gamma) = \text{dnorm}(y, \boldsymbol{\beta}^T \mathbf{x}, \gamma^{-1/2})$ where γ represents the variance. In the Poisson model $p(y | \boldsymbol{\beta}^T \mathbf{x}) = \text{dpois}(\exp\{\boldsymbol{\beta}^T \mathbf{x}\})$, and there is no γ parameter.

```

#### -----
#### EXAMPLE on Linear Mixed Models, HOFF's book - Sect 11.1, 11.3
#### -----
#### -----
#### FUNZIONI che mi servono dopo:
#### Log-density of the multivariate normal distribution
ldmvnorm = function(X, mu, Sigma, iSigma=solve(Sigma), dSigma=det(Sigma)){
  Y = t(t(X)-mu)
  sum(diag(-.5*t(Y)%%Y%%iSigma)) - .5*(prod(dim(X))*log(2*pi) + dim(X)[1]*log(dSigma))
}

### sample from the multivariate normal distribution
rmvnorm = function(n, mu, Sigma){
  p = length(mu)
  res = matrix(0, nrow=n, ncol=p)
  if(n>0 & p>0){
    E = matrix(rnorm(n*p), n, p)
    res = t(t(E)%%chol(Sigma)) + c(mu)
  }
  res
}

### sample from the Wishart distribution
rwish = function(n, nu0, S0){
  S0 = chol(S0)
  S = array(dim=c(dim(S0), n))
  for(i in 1:n){
    Z = matrix(rnorm(nu0*dim(S0)[1]), nu0, dim(S0)[1]) %% S0
    S[,i] = t(Z)%%Z
  }
  S[,1:n]
}

#### -----
#### DATA NESTED in GROUPS (hierarchy), LMM (Linear mixed-effects model)
#### -----
# DATA:
# mathscore dataset of 10-th grade children from 100 different Large urban public high schools
# the data concerns 1993 students
Y = read.table("school_mathscore_01.txt") # This dataset does NOT contain covariate SES

### average values per group of the mathscore
m = length(unique(Y[,1]))
n <- sv <- ybar <- rep(NA, m)
for(j in 1:m) {
  ybar[j] = mean(Y[Y[,1]==j,2]) # mean(Y[[j]])
  sv[j] = var(Y[Y[,1]==j,2])
  n[j] = sum(Y[,1]==j)
}

#### -----
#### LOAD the dataset with covariate SES
#### (it is an indicator of the student's family socioeconomic status)
#### -----
odat = read.table("school_mathscore_02.txt") # This dataset IS INCLUDING covariate values
head(odat)

## sch_id sch_freelunch stu_ses stu_mathscore
## 1 1011 6 -0.0599483 52.11
## 2 1011 6 1.0516522 57.65
## 3 1011 6 -0.8635150 66.44
## 4 1011 6 -0.7965511 44.68
## 5 1011 6 -1.6135105 40.57
## 6 1011 6 -1.1581561 35.04

group = odat$sch_id # numero identificativo (id) della scuola
unique(group)

## [1] 1011 1031 1033 1301 1302 1311 1312 1342 1362 1371 1621 1631 1642 1651 1661
## [16] 1662 1891 1892 1901 1902 1911 1912 1913 2201 2202 2391 2401 2411 2412 2511
## [31] 2522 2541 2542 2551 2552 2562 2563 2671 2683 2781 2782 2791 2792 2801 2802
## [46] 2811 2812 2821 2822 2831 2832 2841 2842 2851 2852 2861 2862 2871 2872 2881
## [61] 2891 3101 3102 3111 3112 3121 3122 3261 3262 3271 3272 3282 3381 3382 3391
## [76] 3392 3401 3402 3411 3412 3421 3432 3441 3451 3452 3561 3562 3571 3631 3632
## [91] 3641 3642 3662 3672 3681 3682 3691 3692 3701 3702

```

```

### -----
### COVARIATES
### -----
# X is a List of m=100 matrices: each matrix has a number of
# rows = number of students in the school (included in the sample)
# each matrix has a number of columns = 2: the first column
# contains 1's, the second contains the value of covariate
# ses (CENTRED wrt the group mean) of each student in the school

X <- list();
for(j in 1:m) {
  xj      <- odat$stu_ses[Y[,1]==j] #Y[,1] è la prima colonna, contiene l'indice della scuola
  xj      <- (xj-mean(xj))
  X[[j]] <- cbind(rep(1,n[j]), xj)
}

ses_cen = odat$stu_ses
summary(ses_cen)

```

```

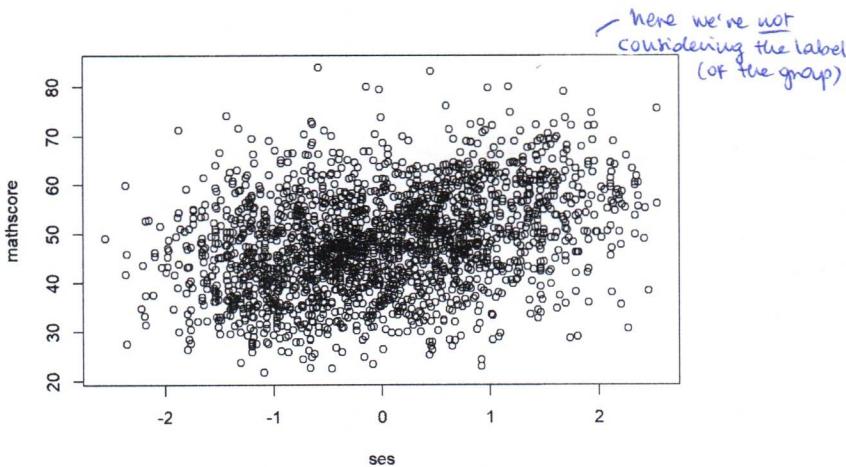
##   Min. 1st Qu. Median Mean 3rd Qu. Max.
## -2.55100 -0.75637 -0.04656 0.00000 0.71683 2.53825

```

```

x11()
### Scatterplot of the data (mathscores and SESS)
plot(odata$stu_ses,Y[,2], xlab='ses', ylab='mathscore')

```



```

### Least Squares Estimates within each group (j=1,...,100)
### perchè potrebbe essere che La relazione tra ses e voto dipenda dalla scuola
S2.LS <- BETA.LS <- NULL
for(j in 1:m) {
  fit     <- lm(Y[,1]==j,2) ~ -1 + X[[j]]) # no intercetta perchè è già nelle matrici X
  BETA.LS <- rbind(BETA.LS, c(fit$coef))
  S2.LS  <- c(S2.LS, summary(fit)$sigma^2)
}

x11()
par(mar=c(2.75,2.75,.5,.5), mgp=c(1.7,.7,0))
par(mfrow=c(1,3))
# LEFT panel: 100 different LS regression Lines AND the average of these lines
# (l'intercetta è la media di tutte le 100 intercette ai min quadrati, e analogamente la slope)
plot(range(ses_cen), range(Y[,2]), type="n", xlab="SES", ylab="math score")
for(j in 1:m){
  abline(BETA.LS[j,1], BETA.LS[j,2], col="gray")
}

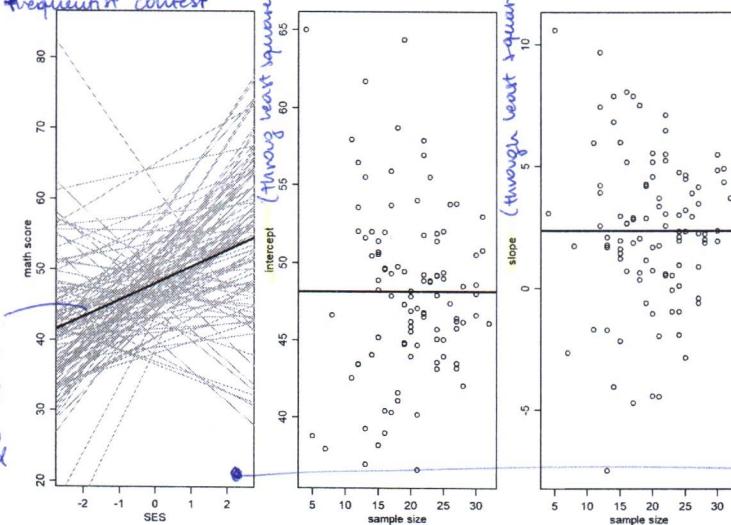
BETA.MLS <- apply(BETA.LS, 2, mean)
abline(BETA.MLS[1], BETA.MLS[2], lwd=2)

# MIDDLE panel
plot(n, BETA.LS[,1], xlab="sample size", ylab="intercept")
# intercepts of the different regression lines vs group (school) sizes
abline(h= BETA.MLS[1], col="black", lwd=2)

# RIGHT panel
plot(n, BETA.LS[,2], xlab="sample size", ylab="slope")
# slopes of the different regression lines vs group (school) sizes
abline(h=BETA.MLS[2], col="black", lwd=2)

```

100 estimated regression lines in the frequentist context



```
## BLACK Lines represent AVERAGE values (average line, average intercept, average slope).
## Schools with the highest sample sizes have regression coefficients that are generally
## close to the average, whereas schools with extreme coefficients are generally those with
## low sample sizes.
## Bayesian solution: stabilize the estimates for small sample size schools by SHARING
## INFORMATION ACROSS GROUPS, using a hierarchical model. Many regression Lines show a positive
## slope, pointing out that, for many groups, as SES increases, maths score increases as well.
## However, there are more than 15 schools with negative slope!
sum(BETA.LS[,2]<0)
```

```
## [1] 19
```

```
### -----
### HIERARCHICAL REGRESSION MODEL - LMM LINEAR MIXED effects MODEL
###

p = dim(X[[1]])[2]
# mu0, the prior expectation of theta is fixed equal to
# the average of the corresponding (frequentist)
# regression parameters
# the matrix Lambda0 is the covariate matrix of these estimates

theta <- mu0 <- apply(BETA.LS, 2, mean)
nu0 <- 1
s2 <- s20 <- mean(S2.LS)
eta0 <- p+2 # così la prior per la matrice Sigma è diffusa
L0 = matrix(nrow=2, ncol=2)
L0[1,1] = cov(BETA.LS)[1,1]
L0[1,2] = cov(BETA.LS)[1,2]
L0[2,1] = cov(BETA.LS)[2,1]
L0[2,2] = cov(BETA.LS)[2,2]

### Inizializzazione della MC (Gibbs Sampler)
Sigma <- S0 <- L0 # <- as.matrix(cov(BETA.LS))
BETA <- BETA.LS
THETA.b <- S2.b <- NULL
iL0 <- solve(L0)
iSigma <- solve(Sigma)
Sigma.ps <- matrix(0,p,p)
SIGMA.PS <- NULL
BETA.ps <- BETA*0
BETA.pp <- NULL
set.seed(1)
mu0[2] + c(-1.96,1.96) * sqrt(L0[2,2]) # prior IC per theta_2 (La slope media)
```

```
## [1] -3.855309 8.602375
```

Some lines have positive slope, others have negative. The case of negative slopes is so much counter-intuitive. Indeed this happens in schools where we have a small number of elements (\rightarrow the frequentist approach is not so strong with small number of elements)



```

### PARAMETERS: beta_1,...,beta_m, theta, SIGMA, sigma^2 (m=100)
### dim(beta_1)=dim(theta)=2, SIGMA 2-by-2 matrix
# sigma ^2 = variance of each response variable (we assume it constant to simplify)
# La sua prior è una inv-gamma(nu0/2,sigma_0^2 nu0/2)
# con nu0=1 e sigma_0^2= media delle varianze campionarie nei gruppi
# beta_1,...,beta,m /theta,SIGMA iid N_2(theta,SIGMA)
# theta ~ N_2(mu0,L0), con mu_0 = vettore delle medie nei gruppi delle stime ai minimi quadrati
# delle intercette e delle slope, mentre L0 = matrice delle covarianze empiriche
# (nei gruppi) delle stime ai minimi quadrati delle intercette e delle slope
# SIGMA ~ inv-Wishart(eta0,5_0^{-1}) con eta_0=4, così E(SIGMA)=1/(eta_0-p-1)* S_0 = matrice delle
# covarianze empiriche delle stime ai minimi quadrati

### Gibbs Sampler cycle
for(s in 1:10000) { # 10000 è un po' troppo, ci mette 5 minuti
  # update beta_j
  for(j in 1:m){
    Vj      <- solve(iSigma + t(X[[j]])%*%X[[j]]/s2)
    Ej      <- Vj%*%(iSigma%*%theta + t(X[[j]])%*%Y[,1]==j,2)/s2)
    BETA[j,] <- rmvnorm(1,Ej,Vj)
  }

  # update theta
  Lm      <- solve(iL0 + m*iSigma)
  mum    <- Lm%*%(iL0%*%mu0 + iSigma%*%apply(BETA,2,sum))
  theta <- t(rmvnorm(1,mum,Lm))

  # update Sigma
  mtheta <- matrix(theta,m,p,byrow=TRUE)
  iSigma <- rwish(1, eta0+m, solve(S0+t(BETA-mtheta)%*%(BETA-mtheta)))

  # update s2
  RSS <- 0
  for(j in 1:m){
    RSS <- RSS + sum((Y[,1]==j,2]-X[[j]]%*%BETA[j,])^2)
  }
  s2 <- 1/rgamma(1,(nu0+sum(n))/2,(nu0*s20+RSS)/2)

  # store results
  if(s%%10==0) {
    cat(s,s2,"\\n")
    S2.b   <- c(S2.b,s2)
    THETA.b <- rbind(THETA.b,t(theta))
    Sigma.ps <- Sigma.ps + solve(iSigma)
    BETA.ps <- BETA.ps + BETA
    SIGMA.PS <- rbind(SIGMA.PS,c(solve(iSigma)))
    BETA.pp <- rbind(BETA.pp,rmvnorm(1,theta,solve(iSigma)))
  }
}

## 10 81.40307
## 20 79.43224
## 30 76.60204
## 40 80.18113
## 50 79.8228
## 60 79.88212
## 70 74.37574
## 80 78.32562
## ..
## 9950 76.44264
## 9960 75.71796
## 9970 77.91109
## 9980 73.92973
## 9990 81.07285
## 10000 80.33422

```

```

## FINE CICLO
## thinning=10, so that the final sample size is 1000

```

```

#save.image("data.f11-3") } la puoi avere già fatto e salvato la run
#Load("data.f11-3")

```

```

### Convergence diagnostics
library(coda)

```

```

effectiveSize(S2.b) # sigma^2

```

```

## var1
## 848.9379

```

```

effectiveSize(THETA.b[,1]) # theta_1

```

```

## var1
## 1000

```

```

effectiveSize(THETA.b[,2]) # theta_2

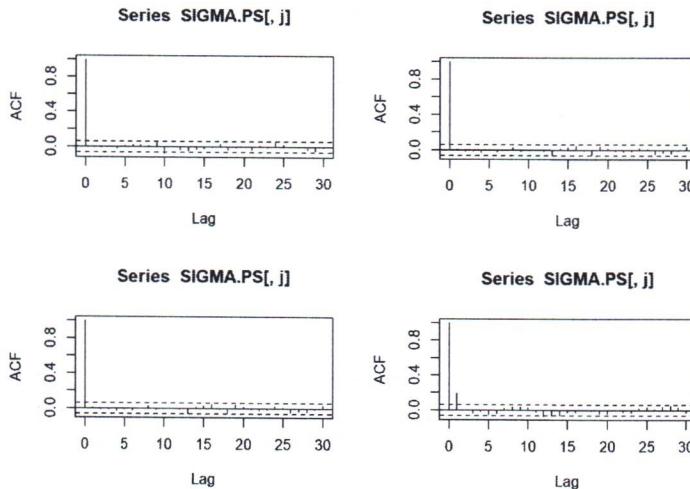
## var1
## 1000

apply(SIGMA.PS,2,effectiveSize) #|\Sigma

## [1] 1000.0000 1000.0000 1000.0000 667.8936

x11()
par(mfrow=c(2,2))
tmp <- NULL
for(j in 1:dim(SIGMA.PS)[2]){
  tmp <- c(tmp,acf(SIGMA.PS[,j])$acf[2])
}

```



```

acf(S2.b)
acf(THETA.b[,1])
acf(THETA.b[,2])

```

```

### -----
# Plot of the marginal posterior of theta2, the population-slope, i.e. the fixed effect
# of covariate ses and the posterior predictive distribution of a NEW school parameter
# (of a school not included in the data), but from the SAME population of schools.
# We sample from the "model" of the beta parameters, when theta, Sigma are the current values
# in the chain. This NEW parameter has been denoted by BETA.pp[,2] in the Gibbs sampler above.
### -----
x11()
par(mar=c(3,3,1,1), mgp=c(1.75,.75,0))
par(mfrow=c(1,2))

plot(density(THETA.b[,2],adj=2), xlim=range(BETA.pp[,2]),
      main="", xlab="slope parameter", ylab="posterior density", lwd=2)
lines(density(BETA.pp[,2],adj=2), col="gray", lwd=2)
legend(-3,1.0, legend=c(expression(theta[2]), expression(tilde(beta)[2])),
       lwd=c(2,2), col=c("black","gray"), bty="n")

quantile(THETA.b[,2], prob=c(.025,.5,.975))

```

```

##      2.5%      50%     97.5%
## 1.832212 2.395010 2.958824

```

```

# 95% posterior CI of the 'average' slope theta_2 is
# very different from the prior CI that is (-3.86, 8.6)

# Posterior probability that theta2 > 0 is very large, but it does
# not indicate that any given within-school slope cannot be negative

# Let us compute the probability that the slope tilde_beta2
# (of a NEW school NOT INCLUDED in the dataset) is negative;
# this value is computed simulating tilde_beta2 from the
# slope population distribution,
# when parameters are the current values of (\thetaeta, \Sigma)
mean(BETA.pp[,2]<0) # this is a small value, but is NOT equal to 0!

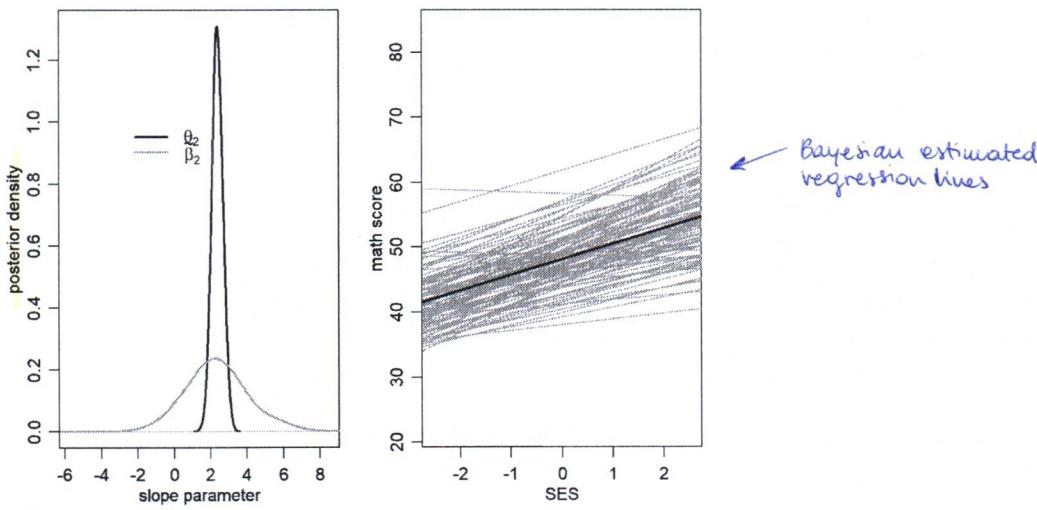
```

```

## [1] 0.072 (small but not zero)

```

```
BETA.PM <- BETA.ps/1000
plot(range(ses_cen), range(Y[,2]), type="n", xlab="SES", ylab="math score")
for(j in 1:m) {
  abline(BETA.PM[j,1], BETA.PM[j,2], col="gray")
}
abline(mean(THETA.b[,1]), mean(THETA.b[,2]), lwd=2)
```

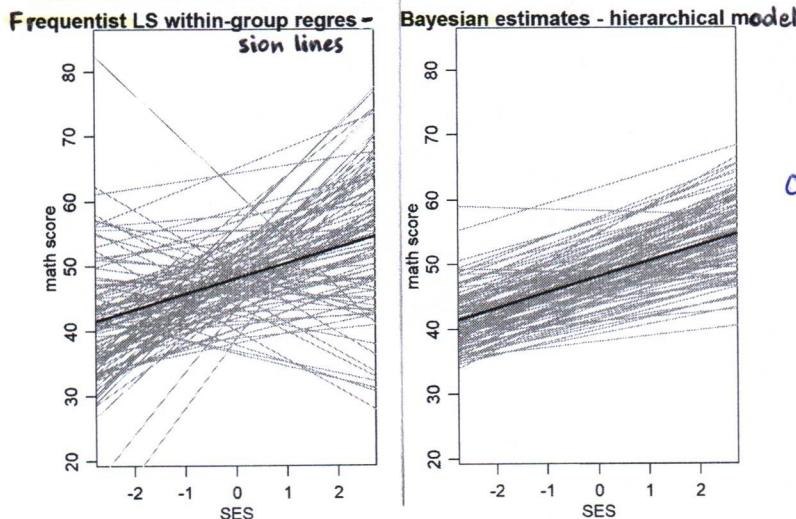


```
### Posterior expectations of the 100 school-specific regression Lines
### i.e. the Lines  $y = E(\beta_{0j})/dati + E(\beta_{1j})/dati \cdot x_{ij}$ 

# Guardate l'analogia figura con le 100 rette di regressione distinte (ai min quadrati)
# RIGHT PLOT: there is a shrinkage of the 'extreme' frequentist regression lines towards the
# across-group average Line (black Line); in fact, the Bayesian estimates we got are convex
# Linear combinations of the (unique) prior mean line and the within-group frequentist estimates.
# Since we SHARED information across groups, hardly any of the slopes are negative!

windows()
par(mar=c(3,3,1,1), mgp=c(1.75,.75,0))
par(mfrow=c(1,2))
plot(range(ses_cen), range(Y[,2]), type="n", xlab="SES", ylab="math score",
     main="Frequentist LS within-group regression lines")
for(j in 1:m) {
  abline(BETA.LS[j,1], BETA.LS[j,2], col="gray")
}

BETA.MLS <- apply(BETA.LS, 2, mean)
abline(BETA.MLS[1], BETA.MLS[2], lwd=2)
plot(range(ses_cen), range(Y[,2]), type="n", xlab="SES", ylab="math score",
     main="Bayesian estimates - hierarchical model")
for(j in 1:m) {
  abline(BETA.PM[j,1], BETA.PM[j,2], col="gray")
}
abline(mean(THETA.b[,1]), mean(THETA.b[,2]), lwd=2)
```



- the lines are estimated through the Bayesian approach (the estimated coefficients for each β_{0j} and β_{1j} are the posterior expectations of those two parameters)

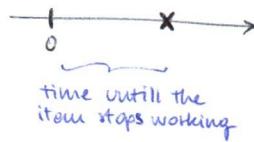
Comment: there are very few lines for which the slope is negative! Moreover the lines seem to be more in agreement (because the Bayesian estimate that we computed were obtained as convex linear combination of the prior mean (which was unique for all the groups) within groups frequentist estimate; we shared information to "help" groups with a small number of values)

```
# LEFT PLOT: 100 frequentist regression Lines
# RIGHT PLOT: there is a shrinkage of the 'extreme' frequentist regression Lines towards the
# across-group average Line (black Line). In fact, the Bayesian estimates we got are convex
# Linear combinations of the (unique) prior mean line and the within-group frequentist estimates.
# Since we SHARED information across groups, hardly any of the slopes are negative!
```

RELIABILITY / SURVIVAL ANALYSIS

Informally, the reliability of an item is the property of the item to work.
What type of data we could have?

- binary → when the item works or not (success / failure)
- time interval when the item work or not
→ date are lifetimes, failure time or time-to-events
(= time when the item stops working)



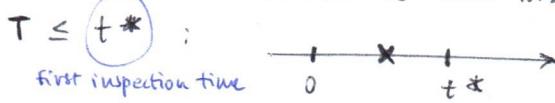
2 features:

1. lifetime / ... are realizations of random variables which are **POSITIVE** and (absolutely) continuous
2. date may be **censored**: date are observed only partially

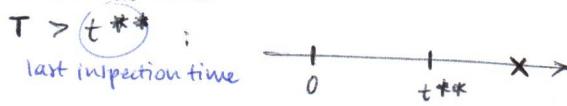
CENSORING:

• LEFT CENSORING

when an item fails before a first inspection ($T = \text{lifetime of the item}$)



• RIGHT CENSORING



: When an item has not failed by the last inspection

• INTERVAL CENSORING:

$$T \in (t^*, t^{**})$$

Focus on **RIGHT-CENSORING**

for some items we know exactly the lifetime,
for some others we know that the event is not occurred up to a certain time

C := random variable which represents the time of censoring :

T := failuretime / lifetime



We typically observe: $y = \min(T, C)$, $\delta = \begin{cases} 1 & \text{if } T \leq c \\ 0 & \text{if } T > c \end{cases}$

so the data is a vector of non-negative data plus an indicator (0/1)

Example: data $(y_0, 0) \iff T > y_0$

- Hp.:
- Independent censoring: T and C are independent
 - Non-informative censoring:

we assume that the distribution of C does not depend on parameters contained in the distribution of T

$C_i(c) = P(C \leq c)$ does not depend on the parameters entering in the distribution of T

Remark: How can we characterize the distribution of a (generic) random variable? Through the cumulative density function. If, in addition, it's an absolutely continuous random variable then we can characterize it through its density. If we have a positive absolutely continuous random variable we can characterize it through the HAZARD FUNCTION.

HAZARD FUNCTION : $h_T(t) = \lambda_T(t) = \frac{f_T(t)}{1 - F_T(t)}$ $t > 0$

it represents the conditional probability of failure per unit of time

$\lim_{\Delta t \rightarrow 0} \frac{\Pr(t \leq T \leq t + \Delta t | T > t)}{\Delta t} = h_T(t)$

= instantaneous probability that the item will stop functioning (why? because)

$$= \lim_{\Delta t \rightarrow 0} \frac{1}{\Delta t} \frac{\Pr(t \leq T \leq t + \Delta t)}{\Pr(T > t)} = \lim_{\Delta t \rightarrow 0} \frac{1}{\Delta t} \frac{F_T(t + \Delta t) - F_T(t)}{\Pr(T > t)}$$

$$= \lim_{\Delta t \rightarrow 0} \frac{F_T(t) \cdot \Delta t}{\Delta t (1 - F_T(t))} = h_T(t)$$

- $F_T \longrightarrow h_T$
- $h_T \longrightarrow F_T(t) = 1 - e^{-\int_0^t h_T(u) du}$
- $\Pr(T > t+s | T > t) = e^{-\int_t^{t+s} h_T(u) du}$

Suppose that: $t \mapsto h_T(t)$ is increasing $\iff \Pr(T > t+s | T > t) \leq \Pr(T > s)$

IFR (Increasing failure rate)
 → we expect to see an increasing number of failures
 (This represents items subjected to aging (e.g. humans))

If instead: $t \mapsto h_T(t)$ is decreasing $\iff \Pr(T > t+s | T > t) \geq \Pr(T > s)$

DFR (Decreasing failure rate)
 → we expect less and less failures in a given period of time

LIKELIHOOD FOR RIGHT-CENSORED DATA

Data: (y_i, δ_i) : $y = (y_1, \dots, y_n)$ $y_i = \min(T_i, C_i)$
 $\delta = (\delta_1, \dots, \delta_n)$ $\delta_i = \begin{cases} 1 & \text{if } T_i \leq C_i \\ 0 & \text{if } T_i > C_i \end{cases}$ = if the event has been observed

Let's assume:

$$T_i | \theta \stackrel{iid}{\sim} f_i(t | \theta) \quad y_i : f_i(t | \theta) = f(t | x_i, \theta)$$

s_i = survival function

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

likelihood given the data

i.e.

- If the i -th observation has been observed
- if the i -th obs. has not been observed

the contribution to the likelihood is:

$$f_i(y_i | \theta) \quad ; \quad s_i(y_i | \theta)$$

 (where $s_i = 1 - F_i$)

Another expression for the likelihood

$$\text{Remember that: } h_i(t | \theta) = \frac{f_i(t | \theta)}{s_i(t | \theta)}$$

$$\Rightarrow f_i(t | \theta) = h_i(t | \theta) \cdot s_i(t | \theta)$$

$$\Rightarrow L(\theta | D) = \prod_{i=1}^n \left\{ (h_i(y_i | \theta))^{\delta_i} (s_i(y_i | \theta))^{1-\delta_i} \right\}$$

$$\Rightarrow L(\theta | D) \propto \prod_{i=1}^n \{ (h_i(y_i | \theta))^{\delta_i} \cdot S_i(y_i | \theta) \}$$

proof. (informal proof of *)

Instead of $\Pr(T_i \in dy | \theta) \rightarrow \Pr(T_i = y | \theta)$. — we consider T_i and C_i "notationally" as discrete random variables (just to avoid mathematics)

Contribution of the i -th observation to the likelihood: → the likelihood is the joint distr. of all the data! Let's see what happens with just one datum

$$\begin{aligned} L(\theta | y_i, \delta_i) &= \begin{cases} \Pr(T_i = y_i, \delta_i = 1 | \theta) & \text{if } \delta_i = 1 \\ \Pr(C_i = y_i, T_i > C_i | \theta) & \text{if } \delta_i = 0 \end{cases} \\ &= \begin{cases} \Pr(T_i = y_i, T_i \leq C_i | \theta) & \text{if } \delta_i = 1 \\ \Pr(C_i = y_i, T_i > C_i | \theta) & \text{if } \delta_i = 0 \end{cases} \\ &= \begin{cases} \Pr(T_i = y_i, C_i \geq y_i | \theta) & \text{if } \delta_i = 1 \\ \Pr(C_i = y_i, T_i > y_i | \theta) & \text{if } \delta_i = 0 \end{cases} \\ &\stackrel{T_i \perp\!\!\!\perp C_i}{=} \begin{cases} \Pr(T_i = y_i | \theta) \Pr(C_i \geq y_i | \theta) & \text{if } \delta_i = 1 \\ \Pr(C_i = y_i | \theta) \Pr(T_i \geq y_i | \theta) & \text{if } \delta_i = 0 \end{cases} \\ &\stackrel{C_i \perp\!\!\!\perp \theta}{=} \begin{cases} f_i(y_i | \theta) \cdot \text{count}_1 & \text{if } \delta_i = 1 \\ S_i(y_i | \theta) \cdot \text{count}_2 & \text{if } \delta_i = 0 \end{cases} \end{aligned}$$

because θ is in the distribution of T_i (second H.p.)
(and since we're looking for the likelihood w.r.t. θ the C_i are constant w.r.t. θ)

we proved what is the contribution to the likelihood in case $\delta_i = 0$ or $\delta_i = 1$

BAYESIAN MODELS FOR RIGHT-CENSORED DATA

Data: $D = (y, \delta)$, $y = (y_1, \dots, y_n)$, $\delta = (\delta_1, \dots, \delta_n)$

$$\begin{aligned} 1. \quad T_1, \dots, T_n | \theta &\stackrel{\text{iid}}{\sim} \mathcal{E}(\theta) \quad \theta > 0 \\ \theta &\sim \text{gamma}(\alpha, \beta) \end{aligned}$$

($\mathcal{E}(\theta)$ is constant (nor IFR or DFR))

$$\begin{aligned} L(\theta | D) &\propto \prod_{i=1}^n (f_i(y_i | \theta))^{\delta_i} (S_i(y_i | \theta))^{1-\delta_i} \\ &= \prod_{i=1}^n (\theta e^{-\theta y_i})^{\delta_i} (e^{-\theta y_i})^{1-\delta_i} \\ &= \theta^{\sum_{i=1}^n \delta_i} \prod_{i=1}^n (e^{-\theta y_i})^{\delta_i + 1 - \delta_i} \\ &= \theta^{\sum_{i=1}^n \delta_i} e^{-\theta \sum_{i=1}^n y_i} \underline{\mathbb{I}_{(0, \infty)}(\theta)} \end{aligned}$$

n_u
(number of uncensored observations)

$$\begin{aligned} f(y | \theta) &= \theta e^{-\theta y} \underline{\mathbb{I}_{(0, \infty)}(y)} \\ S(y | \theta) &= e^{-\theta y} \underline{\mathbb{I}_{(0, \infty)}(y)} \end{aligned}$$

$$\Rightarrow \pi(\theta | D) \propto \theta^{n_u} e^{-\theta \sum_{i=1}^n y_i} \cdot \theta^{\alpha-1} e^{-\beta \theta} \underline{\mathbb{I}_{(0, \infty)}(\theta)}$$

L likelihood · prior

$$\propto \theta^{\alpha + n_u} e^{-(\beta + \sum_{i=1}^n y_i)} \underline{\mathbb{I}_{(0, \infty)}(\theta)}$$

$$\Rightarrow \theta | D \sim \text{gamma}(\alpha + n_u, \beta + \sum_{i=1}^n y_i)$$

The posterior is still a gamma distribution

$$2. T_i | \alpha, \lambda \stackrel{iid}{\sim} \text{Weibull}(\alpha, \lambda) : f(t | \alpha, \lambda) = \lambda \alpha t^{\alpha-1} e^{-\lambda t^\alpha} \mathbb{1}_{(0, \infty)}(t)$$

$$S(t | \alpha, \lambda) = e^{-\lambda t^\alpha} \mathbb{1}_{(0, \infty)}(t)$$

\Rightarrow If $\alpha = 1 \Rightarrow \text{Weibull}(1, \lambda) = \mathcal{E}(\lambda)$

\Rightarrow If $z \sim \mathcal{E}(\lambda) \Rightarrow T = z^{1/\alpha} : T \sim \text{Weibull}(\alpha, \lambda)$

$$h(t | \alpha, \lambda) = \lambda \alpha t^{\alpha-1} \quad t > 0$$

failure rate

- $\alpha < 1 \Rightarrow h(t | \alpha, \lambda)$ is decreasing
 \Rightarrow Weibull DFR

- $\alpha > 1 \Rightarrow h(t | \alpha, \lambda)$ is increasing
 \Rightarrow Weibull IFR

$$\Rightarrow L(\alpha, \lambda | D) \propto (\lambda \alpha)^n \cdot \prod_{i=1}^n (y_i^{\delta_i})^{\alpha-1} e^{-\lambda \sum_{i=1}^n y_i^\alpha}$$

So, what priors do we assume?

$$\pi(\alpha, \lambda) = \pi(\alpha) \pi(\lambda) : \pi(\lambda) = \text{gamma}$$

$$\pi(\alpha) = U([\alpha_1, \alpha_2])$$

such that $1 \in [\alpha_1, \alpha_2]$:

$$0 \quad \alpha_1 \quad 1 \quad \alpha_2 \quad \rightarrow$$

We want to have $[\alpha_1, 1] \neq 0$
 and $[1, \alpha_2] \neq 0$, so that, a posteriori,
 we'll be able to check:

$$P(\alpha < 1 | D), P(\alpha > 1 | D)$$

(we're going to interpret the parameter α
 as the parameter that decides about the
 DFR/IFR type of behaviour)

$$3. T_i \stackrel{iid}{\sim} \text{lognormal}(\mu, \sigma^2) \quad \text{when } T = e^X, X \sim N(\mu, \sigma^2)$$

(lognormal because the log of T is a normal)

prior $\pi(\mu, \sigma^2)$: (either):

- μ, σ^2 independent: $\mu \sim N(\cdot, \cdot)$, $\sigma^2 \sim \text{inv-gamma}$

- $(\mu, \sigma^2) \sim \text{Normal-inv-gamma}$: $\mu | \sigma^2 \sim N(\cdot, \cdot)$, $\sigma^2 \sim \text{inv-gamma}$

$$4. T_i \stackrel{iid}{\sim} \text{gamma}(\alpha, \lambda)$$

- DFR if $\alpha < 1$
- IFR if $\alpha > 1$

(if $\alpha = 1$ we have the gaussian distribution)

$\left. \right\} \Rightarrow$ as for the Weibull, we can put a prior for α and λ
 in such a way that the marginal prior for α keeps
 maps on both the sizes of 1 (and we see then
 what happens a posterior)

TIME-TO-EVENT REGRESSION

ACCELERATED FAILURE TIME MODELS (AFT)

Regression model on
the log-scale of T_i

$$D = \{(y_i, \delta_i, \underline{x}_i)\}_{i=1,\dots,n} : \quad y_i = \min(T_i, c_i)$$

vector of covariates

$$\delta_i = \begin{cases} 1 & \text{if } T_i \leq c_i \\ 0 & \text{if } T_i > c_i \end{cases} \quad \text{= lifetime has been observed}$$

(advice:)
standardize covariates

$$\underline{x}_i = (x_{i1}, \dots, x_{ip}) \quad p = k+1$$

(k = # covariates
 $x_{i1} = 1$)

Model:

$$\log(T_i) = \underline{x}_i^T \beta + \sigma \varepsilon_i \quad \varepsilon_i \stackrel{\text{iid}}{\sim} F_\varepsilon \quad \mathbb{E}[\varepsilon_i] = 0$$

errors

(known) Fixed distribution on \mathbb{R}

$$\beta = (\beta_1, \dots, \beta_p) \quad \leftarrow \text{vector of regression parameters}$$

β_1 intercept

$\sigma > 0$ scale parameter (or τ : $\sigma := \frac{1}{\sqrt{\tau}}$, $\tau = \frac{1}{\sigma^2}$ precision)

$$\Rightarrow T_i | \beta, \tau \stackrel{\text{iid}}{\sim} \text{AFT}(F_\varepsilon, \beta, \tau | \underline{x}_i) \quad \text{— this notation will reserve all the above informations}$$

$$(\beta, \tau) : \quad \bullet \ N_p(0, \dots) \cdot \underline{\pi(\tau)} \cdot \text{gamma}$$

$\bullet \ \pi(\beta | \tau) \cdot \pi(\tau)$ normal-inv-gamma

Why the name "Accelerated Failure Time" model?

$$\text{Notation: } \log(T) = \underline{x}^T \beta + \sigma \varepsilon \quad \varepsilon \sim F_\varepsilon \quad \mathbb{E}[\varepsilon] = 0$$

$$T = e^{\underline{x}^T \beta} \cdot e^{\sigma \varepsilon}$$

$$W := e^\varepsilon \quad W_\tau := e^{\varepsilon \tau} = (W)^\tau$$

$$F_T(t) = \mathbb{P}(T \leq t) = \mathbb{P}(e^{\underline{x}^T \beta} e^{\sigma \varepsilon} \leq t) = \mathbb{P}(e^{\sigma \varepsilon} \leq t e^{-\underline{x}^T \beta}) = F_{W_\tau}(t e^{-\underline{x}^T \beta})$$

accelerate or decelerate the time-scale

What type of distribution do we have for T ? It depends on ε : (F_ε)

$$1. \quad F_\varepsilon = \phi \quad \varepsilon \sim N(0, 1) \quad \rightarrow \quad \sigma \varepsilon \sim N(0, \sigma^2)$$

it should be completely known (no extra parameters)

$$\Rightarrow \log(T_i) \stackrel{\text{iid}}{\sim} N(\underline{x}_i^T \beta, \sigma^2)$$

$\Rightarrow W$ has log-normal distribution
 W_τ has log-normal distribution

Analyzing this distribution means analyzing the distribution of T (and not $\log(T)$)

2. $\varepsilon \sim \text{logistic-distribution}:$

$$f_{\varepsilon}(u) = \frac{e^u}{(1+e^u)^2} \quad u \in \mathbb{R}$$

$$F_{\varepsilon}(u) = \frac{e^u}{1+e^u}$$

$$W := e^{\varepsilon}, \quad W_0 = e^{\sigma\varepsilon}$$

$$F_{W_0}(t) = \mathbb{P}(e^{\sigma\varepsilon} \leq t) = \mathbb{P}\left(\varepsilon \leq \frac{\log(t)}{\sigma}\right) = \frac{e^{\frac{\log(t)}{\sigma}}}{1+e^{\frac{\log(t)}{\sigma}}} = \frac{t^{1/\sigma}}{1+t^{1/\sigma}}$$

3. $\varepsilon \sim \text{extreme value (or Gumbel distribution)}$

$$\begin{aligned} f_{\varepsilon}(u) &= e^u e^{-e^u} \\ F_{\varepsilon}(u) &= 1 - e^{-e^u} \end{aligned} \quad \left. \begin{array}{l} \text{standard} \\ \text{Gumbel} \\ \text{distribution} \end{array} \right.$$

However, we modify it and consider instead:

$$F_{\varepsilon}(u) = 1 - e^{-\log(2)e^u} \quad \leftarrow \text{in this way } 0 \text{ is the median of the distribution}$$

$$\Rightarrow F_{\varepsilon}(0) = \frac{1}{2} \quad \stackrel{\text{"proof":}}{=} F_{\varepsilon}(0) = 1 - e^{-\log(2) \cdot 1} = 1 - \frac{1}{2} = \frac{1}{2}$$

$$W_0 = e^{\sigma\varepsilon}$$

$$F_{W_0}(t) = \mathbb{P}(e^{\sigma\varepsilon} \leq t) = \mathbb{P}\left(\varepsilon \leq \frac{\log(t)}{\sigma}\right) = \underbrace{1 - e^{-\log(2)e^{\frac{\log(t)}{\sigma}}}}_{\substack{= 1 - e^{-\log(2)(t)^{1/\sigma}} \\ t > 0}} \quad \text{Weibull } (\alpha = \frac{1}{\sigma}, \lambda = \log(2))$$

Instead of finding only the distribution of W_0 , let's find the distribution of $T = e^{x^T \beta} e^{\sigma\varepsilon}$:

$$\begin{aligned} F_T(t) &= \mathbb{P}(e^{x^T \beta} e^{\sigma\varepsilon} \leq t) = \mathbb{P}(e^{\sigma\varepsilon} \leq t e^{-x^T \beta}) \\ &\stackrel{!}{=} \mathbb{P}(\sigma\varepsilon \leq \log(t) - x^T \beta) = \mathbb{P}\left(\varepsilon \leq \frac{\log(t)}{\sigma} - \frac{x^T \beta}{\sigma}\right) \\ &= F_{\varepsilon}\left(\frac{\log(t)}{\sigma} - \frac{x^T \beta}{\sigma}\right) = 1 - \exp\left\{-\log(2) e^{\frac{\log(t)}{\sigma}} e^{-\frac{x^T \beta}{\sigma}}\right\} = 1 - \exp\left\{-\log(2) e^{\frac{-x^T \beta}{\sigma}} \cdot (t)^{1/\sigma}\right\} \\ &\quad \text{Weibull } (\alpha = \frac{1}{\sigma}, \lambda = \log(2) e^{-\frac{x^T \beta}{\sigma}}) \end{aligned}$$

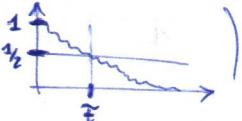
• last important ingredient: MEDIAN FUNCTIONAL OF T

In reliability / survival analysis is often needed a functional of T.

In particular the MEDIAN FUNCTIONAL of T.

$$\tilde{t} \in \mathbb{R}: \quad F_T(\tilde{t}) = \frac{1}{2} \iff S_T(\tilde{t}) = \frac{1}{2}$$

$S_T(\cdot)$ is decreasing. It's of interest to understand when $S_T(\cdot)$ will go from 1 to $\frac{1}{2}$:



\tilde{t} for AFT models:

$$T = e^{x^T \beta} e^{\sigma\varepsilon} \rightarrow \text{we want to know what is } \tilde{t}: \quad F_T(\tilde{t}) = \frac{1}{2}$$

If T has the previous expression, then the distribution of T can be recovered through the distribution of ε :

$$F_T(t) = F_{\varepsilon}\left(\frac{\log(t) - x^T \beta}{\sigma}\right) = \frac{1}{2}$$

If 0 is the median of the distr. of ε , i.e. $F_{\varepsilon}(0) = \frac{1}{2}$

$$\Rightarrow \tilde{t}: \quad \frac{\log(\tilde{t}) - x^T \beta}{\sigma} = 0 \quad \Rightarrow \quad \boxed{\tilde{t} = e^{x^T \beta}}$$

We will monitor what happens with this $g(\beta) = e^{x^T \beta}$
(in a bayesian approach, $g(\cdot)$ has a prior)

 So e^{β_2} is the effect of stage 2 w.r.t. stage 1 over the median (is the ratio of two medians of patients which have exactly the same covariates but one is stage 2 and the other is stage 1).

Similarly we'll monitor e^{β_3} : e^{β_3} = relative median of two patients for which the covariates are all equal but the patient at the numerator has tumor at stage 3 while patient at the denominator has stage 1.

Similarly with e^{β_4} .

What is e^{β_5} ? It's the relative median of two patients which have ^{all} same covariates but the age. (Remember that we standardize covariates) ~~so the difference between their age will be 0 or 1~~

At the numerator we'll have a patient which has:

$$\text{age-num} = (\text{age-den}) + (1 \text{ standard deviation})$$

Similarly for e^{β_6} .

We'll monitor if all these ratios will be ≈ 1 .
(From that it'll be clear ~~which~~ which change in the covariates has effect in increasing/decreasing the median)

~~$$\text{age} = \text{age at birth} + 1 \text{ std deviation of age}$$~~

Survival probability at 5 months for new patients:

$$3. \quad \text{IP}(\tau_{\text{new}} > 5 \text{ months} | x_{\text{new}}, \beta, \tau) = \\ = \exp \left\{ -\log(2) e^{(\log(S) - x_{\text{new}} \beta) / \tau} \right\}$$

Data:
 some date are fully observed: t_1, \dots, t_{n^*}
 some instead are right censored
 $\wedge T_{n^*+1}, T_{n^*+2}, \dots, T_n \rightarrow C_{n^*+j}$
 we only know that:
 (small c_i , we already know the value)

Jags has to write down the likelihood:

$$\prod_{i=1}^n \left\{ (f_i(t|\theta))^{\delta_i} (S_i(t|\theta))^{1-\delta_i} \right\}$$

as bayesians, we'll treat these values as missing data *

and so he needs to know how to treat *.
 (We must tell him how to deal with 'em)

```

### -----
### -----
### AFT models
### -----
# We read data and invoke jags from within R
rm(list=ls())
library(rjags)

```

```

in_data = read.table("data_larynx_cancer.txt", header=TRUE)
head(in_data)

## stage t age Yr cens
## 1 1 0.6 77 76 0.0
## 2 1 1.3 53 71 0.0
## 3 1 2.4 45 71 0.0
## 4 1 NA 57 78 2.5
## 5 1 3.2 58 74 0.0
## 6 1 NA 51 77 3.2

```

Stage of the tumor when diagnosed (1,2,3,4)
 age of the patient when diagnosed
 year of diagnosis
 the patients died in the first months and we observed it
 for this patient we only know that his survival time was greater than 2.5 but at that time the patient left the study so we don't know anything else

```

# Data on 90 male patients
# Lifetime= time in months from diagnosis to death or censoring
# stage t age Yr cens
#1 1 0.6 77 76 0.0
#2 1 1.3 53 71 0.0
#3 1 2.4 45 71 0.0
#4 1 NA 57 78 2.5 right censored
# Subjects 1,2,3 have observed Lifetime, while
# subject 4 is censored, i.e. T_i>2.5

## data on 90 male patients with cancer of Larynx: 50 survival
## times are observed, 40 are right censored
## survival times(in months)
## COVARIATE: stage of the disease at diagnosis (1,2,3,4) (stage)
## year of diagnosis (Yr)
## age at diagnosis (age)
## Extra 16 rows: covariates values for 16 NEW patients; we will
## compute their predictive distribution
## In the bug file we will standardize Age and Yr
## For people with the same age and year of diagnosis, we are
## interested in the median Lifetimes for individuals
## in Stage 2,3 or 4 relative to Stage 1
## We are also interested in the predictive 5-months survival for 16 NEW patients

attach(in_data)
# HERE in this dataset cens= 0 means that datum has been OBSERVED

# we need to adjust this vector for JAGS!!
tmax      <- round(max(t[!is.na(t)]))+1
cens[!is.na(t)] <- tmax
tmax

```

and so we add 16 rows at the end of the dataset (4 with diagnosis "1", ..., 4 with diagnosis "4", with combinations in age (50/70) and year (71/77))

```

## [1] 9

```

```

cens

```

```

## [1] 9.0 9.0 9.0 2.5 9.0 3.2 9.0 3.3 9.0 9.0 9.0 9.0 9.0 4.5 4.5
## [16] 9.0 5.5 5.9 5.9 9.0 6.1 6.2 9.0 9.0 6.5 6.7 7.0 9.0 7.4 8.1
## [31] 8.1 9.6 10.7 9.0 9.0 9.0 2.2 2.6 3.3 9.0 3.6 9.0 4.3 4.3 5.0
## [46] 9.0 9.0 7.5 7.6 9.3 9.0 9.0 9.0 9.0 9.0 9.0 9.0 9.0 9.0 9.0
## [61] 9.0 9.0 9.0 3.7 4.5 4.8 4.8 9.0 5.0 5.1 9.0 9.0 6.5 9.0 8.0
## [76] 9.3 10.1 9.0 9.0 9.0 9.0 9.0 9.0 9.0 9.0 9.0 2.9 9.0 9.0 4.3
## [91] 0.0 0.0 0.0 0.0 0.0 0.0 0.0 0.0 0.0 0.0 0.0 0.0 0.0 0.0 0.0
## [106] 0.0

```

```

## Censored data must be recorded as NA, not as the value of censoring Limit.
# When explicitly initializing the chains, the censored values of
# the data must be explicitly initialized (to values above the censoring limits)!
# SEE http://doingbayesiandataanalysis.blogspot.com/2012/01/complete-example-of-right-censoring-in.html
# or Julian Stander, Luciana Dalla Valle & Mario Cortina-Borja (2018).
# A Bayesian Survival Analysis of a Historical Dataset: How Long
# Do Popes Live?, The American Statistician, 72:4, 368-375, DOI: 10.1080/00031305.2017.1328374
head(cens)

```

```

## [1] 9.0 9.0 9.0 2.5 9.0 3.2

```

```

dat      = list(cens=cens, t=t, stage=stage, Yr=Yr, age=age)
inits   = function() {list(beta=c(2,-0.14,-0.5,-1.5,-0.2,0), alpha=1)}
modelAFT = jags.model("larynx-AFT.jags", data=dat, n.chains=3)

```

```

## Compiling model graph
##   Resolving undeclared variables
##   Allocating nodes
## Graph information:
##   Observed stochastic nodes: 50
##   Unobserved stochastic nodes: 63
##   Total graph size: 1444
##
## Initializing model

update(modelAFT,50000)

# We monitor beta, sigma, rm, prob, mu and S
variable.names = c("beta", "alpha", "rm", "prob", "mu", "S") # monitoring
n.iter      = 50000
thin        = 10
outputAFT   = coda.samples(model=modelAFT,variable.names=variable.names,n.iter=n.iter,thin=thin)

#save(outputAFT,file='AFT_output.Rdata') # we save the chain
#load('AFT_output.Rdata')

library(coda)
data.out <- as.matrix(outputAFT)
data.out <- data.frame(data.out)
attach(data.out)
n.chain <- dim(data.out)[1]
n.chain

## [1] 15000

summary(data.out)

##      S.91.          S.92.          S.93.          S.94.
##  Min. :0.01547  Min. :0.0000000  Min. :0.09675  Min. :0.00000
##  1st Qu.:0.25487  1st Qu.:0.0007894  1st Qu.:0.39122  1st Qu.:0.00616
##  Median :0.35149  Median :0.0061268  Median :0.46615  Median :0.02379
##  Mean   :0.35839  Mean   :0.0252145  Mean   :0.46581  Mean   :0.04844
##  3rd Qu.:0.45488  3rd Qu.:0.0277426  3rd Qu.:0.53895  3rd Qu.:0.06588
##  Max.  :0.85879  Max.  :0.5559940  Max.  :0.83968  Max.  :0.61978
##      S.95.
##      ..
```

```

names(data.out)

##  [1] "S.91."    "S.92."    "S.93."    "S.94."    "S.95."    "S.96."    "S.97."
##  [8] "S.98."    "S.99."    "S.100."   "S.101."   "S.102."   "S.103."   "S.104."
## [15] "S.105."   "S.106."   "alpha"    "beta.1."  "beta.2."  "beta.3."  "beta.4."
## [22] "beta.5."  "beta.6."  "mu.1."   "mu.2."   "mu.3."   "mu.4."   "mu.5."
## [29] "mu.6."   "mu.7."   "mu.8."   "mu.9."   "mu.10."  "mu.11."  "mu.12."
## [36] "mu.13."  "mu.14."  "mu.15."  "mu.16."  "mu.17."  "mu.18."  "mu.19."
## [43] "mu.20."  "mu.21."  "mu.22."  "mu.23."  "mu.24."  "mu.25."  "mu.26."
## [50] "mu.27."  "mu.28."  "mu.29."  "mu.30."  "mu.31."  "mu.32."  "mu.33."
## [57] "mu.34."  "mu.35."  "mu.36."  "mu.37."  "mu.38."  "mu.39."  "mu.40."
## [64] "mu.41."  "mu.42."  "mu.43."  "mu.44."  "mu.45."  "mu.46."  "mu.47."
## [71] "mu.48."  "mu.49."  "mu.50."  "mu.51."  "mu.52."  "mu.53."  "mu.54."
## [78] "mu.55."  "mu.56."  "mu.57."  "mu.58."  "mu.59."  "mu.60."  "mu.61."
## [85] "mu.62."  "mu.63."  "mu.64."  "mu.65."  "mu.66."  "mu.67."  "mu.68."
## [92] "mu.69."  "mu.70."  "mu.71."  "mu.72."  "mu.73."  "mu.74."  "mu.75."
## [99] "mu.76."  "mu.77."  "mu.78."  "mu.79."  "mu.80."  "mu.81."  "mu.82."
## [106] "mu.83."  "mu.84."  "mu.85."  "mu.86."  "mu.87."  "mu.88."  "mu.89."
## [113] "mu.90."  "mu.91."  "mu.92."  "mu.93."  "mu.94."  "mu.95."  "mu.96."
## [120] "mu.97."  "mu.98."  "mu.99."  "mu.100." "mu.101." "mu.102." "mu.103."
## [127] "mu.104." "mu.105." "mu.106." "prob.1."  "prob.2."  "prob.3."  "prob.4."
## [134] "prob.5."  "prob.6."  "rm.1."   "rm.2."   "rm.3."   "rm.4."   "rm.5."
## [141] "rm.6."
```

```

#####
### beta parameters
#####
summary(data.out[,18:23])

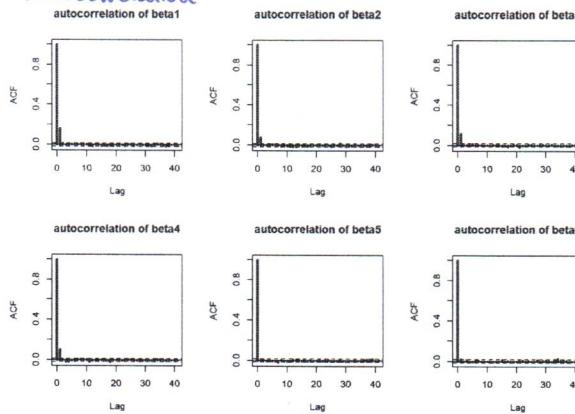
##      beta.1.          beta.2.          beta.3.          beta.4.
##  Min. :-0.1063  Min. :-1.21950  Min. :-1.25649  Min. :-1.7334
##  1st Qu.: 0.7659  1st Qu.:0.07569  1st Qu.:-0.41637  1st Qu.:-0.7544
##  Median : 0.9019  Median :0.14339  Median :-0.24911  Median :-0.5442
##  Mean   : 0.9028  Mean   : 0.15187  Mean   :-0.25521  Mean   :-0.5479
##  3rd Qu.: 1.0366  3rd Qu.: 0.37008  3rd Qu.:-0.08687  3rd Qu.:-0.3437
##  Max.  : 1.7617  Max.  : 2.05414  Max.  : 0.68261  Max.  : 0.7705
##      beta.5.          beta.6.
##  Min. :-0.38149  Min. :-0.9518
##  1st Qu.: 0.04382  1st Qu.:-0.5115
##  Median : 0.11476  Median :-0.4239
##  Mean   : 0.11313  Mean   :-0.4220
##  3rd Qu.: 0.18542  3rd Qu.:-0.3316
##  Max.  : 0.63814  Max.  : 0.2423
```

```

x11()
par(mfrow=c(2,3))
acf(data.out[, 'beta.1.'], lwd=3, col="red3", main="autocorrelation of beta1")
acf(data.out[, 'beta.2.'], lwd=3, col="red3", main="autocorrelation of beta2")
acf(data.out[, 'beta.3.'], lwd=3, col="red3", main="autocorrelation of beta3")
acf(data.out[, 'beta.4.'], lwd=3, col="red3", main="autocorrelation of beta4")
acf(data.out[, 'beta.5.'], lwd=3, col="red3", main="autocorrelation of beta5")
acf(data.out[, 'beta.6.'], lwd=3, col="red3", main="autocorrelation of beta6")

```

Autocorrelation

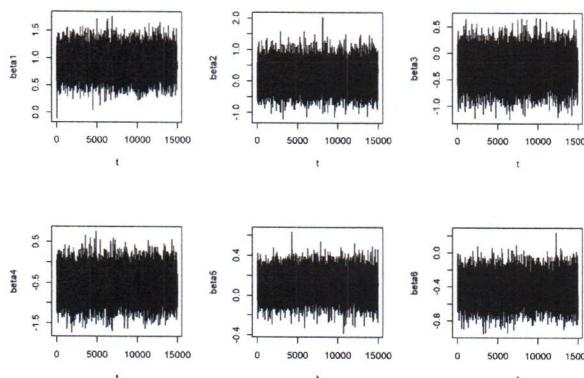


```

x11()
par(mfrow=c(2,3))
plot(ts(data.out[, 'beta.1.']), xlab="t", ylab="beta1")
plot(ts(data.out[, 'beta.2.']), xlab="t", ylab="beta2")
plot(ts(data.out[, 'beta.3.']), xlab="t", ylab="beta3")
plot(ts(data.out[, 'beta.4.']), xlab="t", ylab="beta4")
plot(ts(data.out[, 'beta.5.']), xlab="t", ylab="beta5")
plot(ts(data.out[, 'beta.6.']), xlab="t", ylab="beta6")

```

Traceplot

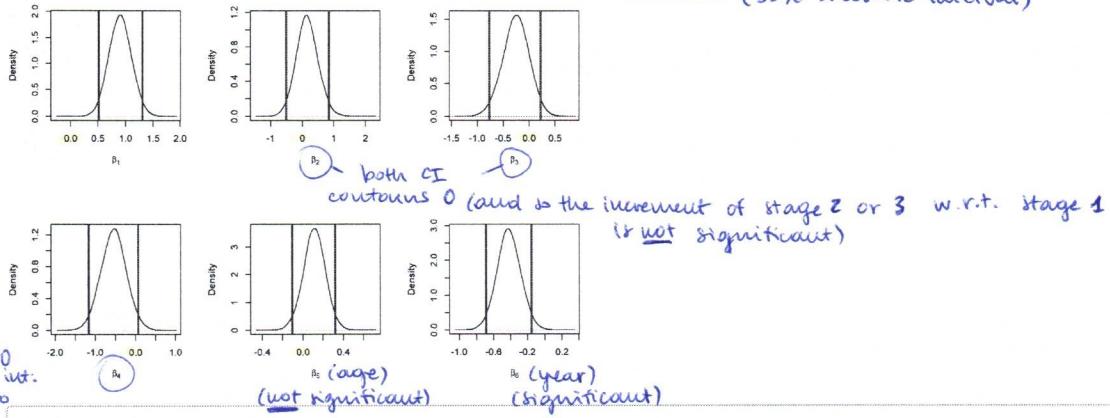


```

# Kernel density plot for marginal posteriors of beta coefficients
x11()
par(mfrow=c(2,3))
plot(density(data.out[, 'beta.1.'], adj=2), xlab=expression(beta[1]), main="")
abline(v=quantile(data.out[, 'beta.1.'], prob=c(.025, .975)), lwd=2, col="red")
plot(density(data.out[, 'beta.2.'], adj=2), xlab=expression(beta[2]), main="")
abline(v=quantile(data.out[, 'beta.2.'], prob=c(.025, .975)), lwd=2, col="red")
plot(density(data.out[, 'beta.3.'], adj=2), xlab=expression(beta[3]), main="")
abline(v=quantile(data.out[, 'beta.3.'], prob=c(.025, .975)), lwd=2, col="red")
plot(density(data.out[, 'beta.4.'], adj=2), xlab=expression(beta[4]), main="")
abline(v=quantile(data.out[, 'beta.4.'], prob=c(.025, .975)), lwd=2, col="red")
plot(density(data.out[, 'beta.5.'], adj=2), xlab=expression(beta[5]), main="")
abline(v=quantile(data.out[, 'beta.5.'], prob=c(.025, .975)), lwd=2, col="red")
plot(density(data.out[, 'beta.6.'], adj=2), xlab=expression(beta[6]), main="")
abline(v=quantile(data.out[, 'beta.6.'], prob=c(.025, .975)), lwd=2, col="red")

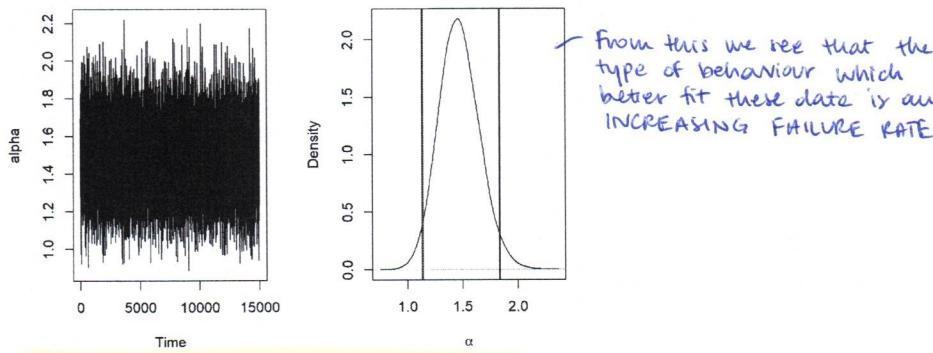
```

First thing: check if 0 is in the posterior credible interval (95% credible interval)



```
###  
### alpha = 1/5  
###  
x11()  
par(mfrow=c(1,2))  
plot(ts(data.out[, 'alpha']), ylab="alpha")  
  
plot(density(data.out[, 'alpha'], adj=2), xlab=expression(alpha), main="Aging Parameter of Weibull")  
abline(v=quantile(data.out[, 'alpha'], prob=c(.025, .975)), lwd=2, col="red")
```

Aging Parameter of Weibull



We monitor the RELATIVE MEDIAN:

```
# rm[2]=exp(beta[2]): relative median of a patient with Stage 2  
# wrt to a patient (same age and Yr) with Stage 1,  
# i.e. ratio of medians of two individuals with same age and Yr,  
# but one with Stage 2 (numerator) and the other (denom) with Stage 1.  
# rm[3]=exp(beta[3]) and rm[4]=exp(beta[4]) measure the same but  
# comparing Stages 3 and 4 versus 1, respectively.  
# rm[5]=exp(beta[5]): relative median of two individuals who are  
# 1 (sample) standard deviation apart in age (=10.8years),  
# but with the same stage and year of diagnosis  
# rm[6]=exp(beta[6]): same but for year of diagnosis (sd=2.34)  
#  
rmchain = cbind(rm_2=data.out[, 'rm.2.'], rm_3=data.out[, 'rm.3.'],  
rm_4 = data.out[, 'rm.4.'], rm_5=data.out[, 'rm.5.'], rm_6=data.out[, 'rm.6.'])  
apply(rmchain, 2, "quantile", prob=c(0.025, 0.5, 0.975))
```

```
## rm_2 rm_3 rm_4 rm_5 rm_6  
## 2.5% 0.5999879 0.4636210 0.3124000 0.8995456 0.5013697  
## 50% 1.1541816 0.7794954 0.5803276 1.1215987 0.6544844  
## 97.5% 2.3267219 1.2524925 1.0700910 1.3819058 0.8581131
```

```
###  
### Posterior probability that beta_i>0  
###  
# beta_1 represents the average effect (give the Age and Yr)  
# of Stage 1 on the Log-survival time;  
# beta_1+beta_2 represents the average effect (give the Age and Yr)  
# of Stage 2 on the Log-survival time; so beta_2 is the effect  
# of Stage 2 wrt Stage 1  
probchain = cbind(data.out[, 'prob.2.'], data.out[, 'prob.3.'], data.out[, 'prob.4.'],  
data.out[, 'prob.5.'], data.out[, 'prob.6.'])  
apply(probchain, 2, mean)
```

```
## [1] 0.67140000 0.15273333 0.03953333 0.85800000 0.00120000
```

```
apply(probchain, 2, sd)
```

```
## [1] 0.46971986 0.35974225 0.19486657 0.34906178 0.03462138
```

```

# beta_3 is NEGATIVE, but with uncertainty
# beta_4 is NEGATIVE, but with less uncertainty
# beta_5, the age effect, is positive

### -----
### Predictions for "new" patients
### -----
## Covariate contenute nelle ultime 16 righe (dalla riga 91 alla riga 106) di data_in
## Stage = 1,2,3,4; Age = 50,70; Yr = 71,77
in_data[91:106,]

```

```

##   stage t age Yr cens
## 91     1 NA 50 71    0
## 92     1 NA 50 77    0
## 93     1 NA 70 71    0
## 94     1 NA 70 77    0
## 95     2 NA 50 71    0
## 96     2 NA 50 77    0
## 97     2 NA 70 71    0
## 98     2 NA 70 77    0
## 99     3 NA 50 71    0
## 100    3 NA 50 77    0
## 101    3 NA 70 71    0
## 102    3 NA 70 77    0
## 103    4 NA 50 71    0
## 104    4 NA 50 77    0
## 105    4 NA 70 71    0
## 106    4 NA 70 77    0

```

```
names(data.out)
```

```

## [1] "S.91."   "S.92."   "S.93."   "S.94."   "S.95."   "S.96."   "S.97."
## [8] "S.98."   "S.99."   "S.100."  "S.101."  "S.102."  "S.103."  "S.104."
## [15] "S.105."  "S.106."  "alpha"   "beta.1." "beta.2." "beta.3." "beta.4."
## [22] "beta.5." "beta.6." "mu.1."   "mu.2."   "mu.3."   "mu.4."   "mu.5."
## [29] "mu.6."   "mu.7."   "mu.8."   "mu.9."   "mu.10."  "mu.11."  "mu.12."
## [36] "mu.13."  "mu.14."  "mu.15."  "mu.16."  "mu.17."  "mu.18."  "mu.19."
## [43] "mu.20."  "mu.21."  "mu.22."  "mu.23."  "mu.24."  "mu.25."  "mu.26."
## [50] "mu.27."  "mu.28."  "mu.29."  "mu.30."  "mu.31."  "mu.32."  "mu.33."
## [57] "mu.34."  "mu.35."  "mu.36."  "mu.37."  "mu.38."  "mu.39."  "mu.40."
## [64] "mu.41."  "mu.42."  "mu.43."  "mu.44."  "mu.45."  "mu.46."  "mu.47."
## [71] "mu.48."  "mu.49."  "mu.50."  "mu.51."  "mu.52."  "mu.53."  "mu.54."
## [78] "mu.55."  "mu.56."  "mu.57."  "mu.58."  "mu.59."  "mu.60."  "mu.61."
## [85] "mu.62."  "mu.63."  "mu.64."  "mu.65."  "mu.66."  "mu.67."  "mu.68."
## [92] "mu.69."  "mu.70."  "mu.71."  "mu.72."  "mu.73."  "mu.74."  "mu.75."
## [99] "mu.76."  "mu.77."  "mu.78."  "mu.79."  "mu.80."  "mu.81."  "mu.82."
## [106] "mu.83."  "mu.84."  "mu.85."  "mu.86."  "mu.87."  "mu.88."  "mu.89."
## [113] "mu.90."  "mu.91."  "mu.92."  "mu.93."  "mu.94."  "mu.95."  "mu.96."
## [120] "mu.97."  "mu.98."  "mu.99."  "mu.100." "mu.101." "mu.102." "mu.103."
## [127] "mu.104." "mu.105." "mu.106." "prob.1." "prob.2." "prob.3." "prob.4."
## [134] "prob.5." "prob.6." "rm.1."  "rm.2."  "rm.3."  "rm.4."  "rm.5."
## [141] "rm.6."

```

```

# Posterior credibility regions for predicted survivals
# (last 16 rows of the dataset)
predsurv <- apply(data.out[,1:16],2,"quantile",prob=c(0.025,0.5,0.975))
print(predsurv)

```

```

##      S.91.     S.92.     S.93.     S.94.     S.95.     S.96.
## 2.5% 0.1116095 1.509764e-06 0.2611363 0.0001370411 0.0921594 2.314117e-05
## 50% 0.3514944 6.126802e-03 0.4661521 0.0237927131 0.4315190 1.652908e-02
## 97.5% 0.6430849 1.646696e-01 0.6711099 0.2341249394 0.7741743 2.561660e-01
##      S.97.     S.98.     S.99.     S.100.    S.101.    S.102.
## 2.5% 0.2200715 0.001100828 0.04467635 1.632779e-08 0.1392830 3.321398e-06
## 50% 0.5402983 0.048509125 0.222723088 6.369331e-04 0.3319497 4.538722e-03
## 97.5% 0.8039518 0.322121285 0.50619959 6.181283e-02 0.5647889 1.065942e-01
##      S.103.    S.104.    S.105.    S.106.
## 2.5% 0.001136063 2.244360e-13 0.01801223 1.469499e-08
## 50% 0.098277538 1.165394e-05 0.18285083 2.501591e-04
## 97.5% 0.478377158 1.633162e-02 0.51235150 2.642605e-02

```

for each patient (new)
there are the credibility
intervals of the survival
probability of 5 months

```

#### -----
#### ----- JAGS code for AFT model for the larynx cancer dataset
#### -----
#### -----
# AFT model for data set on 90 males with larynx cancer.
# Model is a log-linear type of specification (AFT).
# lifetime= time in months from diagnosis to death or censoring
# Covariates are:
# - Stage of the disease at diagnosis time (one through four) [coded using dummies]
# - Age of patient at diagnosis
# - Yr, the year of diagnosis
#
# # Age and Yr will be standardized (i.e. their means and sds are 0 and 1, respectively).
#
# The data also include cens[i], defined as 0 if not censored, tmax if censored
# Note: the last 16 lines of the data file contain points
#      at which we wish to make predictions for some combinations
#      of Stage, Age and Yr, i.e. they are not actual data!
#
# The model uses a linear predictor function
#   mu[i] = beta1 + beta2 * S2[i] + beta3 * S3[i] + beta4 * S4[i]
#   + beta5 * Age[i] + beta6 * Yr[i]
#
# We will use the specification T[i] ~ Weibull(alpha,lambda[i]), with
# lambda[i] = log(2) * exp(-mu[i] * alpha)
#
# Weibull distribution has density parametrized in JAGS as
#   lambda * alpha * t^{alpha-1} * exp(-lambda * t^alpha)
# and CDF = 1 - exp(-lambda * t^alpha).
# The median is exp(mu[i]).
#
# The resulting distribution for log(T[i]) has CDF
#   1 - exp(-log(2) * e^{alpha*(t-mu[i])}) = 1 - 2**[-exp{alpha*(t-mu[i])}]
# so the median is mu[i]. Thus:
#
# P(T[i] > t) = P(log(T[i]) > log(t)) which we can use as survival function.
#### -----
#### The MODEL in JAGS
#### -----
model {
for(i in 1:106){
  sAge[i] <- (age[i]-mean(age[ ]))/sd(age[ ]) # standardization
  sYr[i] <- (Yr[i]-mean(Yr[ ]))/sd(Yr[ ]) # standardization
  censored[i] ~ dinterval(t[i], cens[i]) # this is how JAGS understands censoring
  t[i] ~ dweib(alpha, lambda[i])
}
# The indicator censored takes the value 1 for censored data, 0 for the others.
# The function dinterval(s,c) takes the value 0 if s <= c and 1 if s > c.
# When i = 4 (censored data, t[4]=NA), the requirement
# censored[4] ~ dinterval(t[4], censoring_limits[4])
# becomes 1 ~ dinterval(NA, 2.5), which forces the unknown t[4]
# to be simulated from the Weibull subject to the constraint that t[4] > 2.5.
# When i=1 (observed data), the requirement
# censored[i] ~ dinterval(t[i], cens[i]) becomes
# 0 ~ dinterval(t[i], cens[1]).
# In general, this requirement would be automatically satisfied
# for any value t[1] (or t[i] more generally) greater than or equal
# to all the noncensored survival times if cens[1] (and cens of all
# observed data) is LARGER than all observed data, so the
# actual value chosen for any cens[i] corresponding to
# observed data is unimportant provided that it is sufficiently
# large:
# in the R script, we have set
# tmax <- round(max(t[is.na(t)]))+1
# cens[!is.na(t)] <- tmax, so that cens[1]=9
#
# lambda[i] <- log(2)*exp(-mu[i]*alpha)
# re-parametrized Weibull, median centered
# Covariate Stage is categorical (3 dummy variables)
  mu[i] <- beta[1] + beta[2]*equals(stage[i],2) + beta[3]*equals(stage[i],3) +
    beta[4]*equals(stage[i],4) + beta[5]*sAge[i] + beta[6]*sYr[i]
}
# 5 month survival probabilities corresponding to predictors in augmented data
for(i in 91:106) {
  S[i] <- exp(-log(2)*exp(( log(5) - mu[i])*alpha))
# And estimated medians in original scale
  med[i] <- exp(mu[i])
}
for(i in 1:6){
  beta[i] ~ dnorm(0,0.0001) # 0.001
  rm[i] <- exp(beta[i])
  prob[i] <- step(beta[i])
}
# Relative medians for each variable.
# Note: for covariates x1 and x2 the relative medians are
# exp(x1*beta)/exp(x2*beta)= exp((x1-x2)*beta), so we get
# relative medians by taking x1=(0,1,0,0,0,0), x2=(0,0,0,0,0,0), ...
#
# prob[i] is the posterior probability that beta[i] > 0
}
alpha ~ dunif(0,10) # alpha= 1/sigma
}

```

the " ~ " does not stand for distribution, it represents a condition to be respected

W is usually symmetric

W can be standardized: $\tilde{W} = [\tilde{w}_{ij}]$

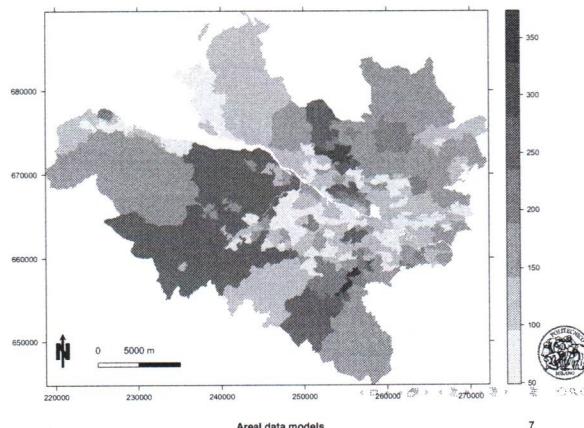
$$\tilde{w}_{ij} = w_{ij}/w_{i+}, \quad w_{i+} := \sum_j w_{ij}$$

(the row sum is always 1, but \tilde{W} is no more symmetric)

w_{ij} can be viewed as a weight: more weight will be associated with j 's closer (in some sense) to i than those farther away from i



Median property prices in 270 areal units



Measures of spatial association

$$\text{Moran's } I = \frac{n \sum_i \sum_j w_{ij} (Y_i - \bar{Y})(Y_j - \bar{Y})}{(\sum_{i \neq j} w_{ij}) \sum_i (Y_i - \bar{Y})^2}$$

analogue of lagged autocorrelation in time series

$$\text{Geary's } C = \frac{(n-1) \sum_i \sum_j w_{ij} (Y_i - \bar{Y})^2}{2(\sum_{i \neq j} w_{ij}) \sum_i (Y_i - \bar{Y})^2}$$

For large n , under the null H_0 that Y_i 's are iid:

$$I \sim \mathcal{N}\left(-\frac{1}{n-1}, \text{Var}(I)\right)$$

$$C \sim \mathcal{N}(1, \text{Var}(C))$$

`moran.test` and `geary.test` in `spdep`



Full conditional distributions

Which joint distribution (likelihood) should we choose for

$$\mathbf{Y} = (Y_1, \dots, Y_n)$$

If we give the full conditionals $\mathcal{L}(Y_i | Y_{-i})$, i.e. $p(y_i | y_{-i})$, for all $i = 1, \dots, n$, do we determine the joint law $p(y_1, \dots, y_n)$ of \mathbf{Y} ?

YES (under consistency conditions of the full-conditionals), since $p(y_1, \dots, y_n)$ is determined starting from the product of all $p(y_i | y_{-i})$, but the joint distribution could be improper

we assume the full conditionals to be gaussian, so this property is guaranteed

Remark: if n is large, we do not seek to write down the joint distribution of \mathbf{Y} ; we prefer to work and model exclusively with the n full-conditionals, since they may represent the *local* behaviour of each Y_i



Second famous model (the first was AFT)

COX-MODEL or PROPORTIONAL HAZARD MODEL (PH model)

27/11/20

Similar context as before (right censored data and covariates). We're going to model T (the distr. of T) modelling the hazard function.

T :
$$h_T(t|x, \beta) = e^{x^T \beta} \cdot h_0(t)$$

Hazard function of T baseline hazard function
 f_0, S_0 (corresponding density and survival function of h_0)

$x^T \beta$ increases $\Rightarrow h(t|x, \beta)$ increases as well

DO NOT include the ~~intercept~~ intercept in β in the proportional hazard model

$$S(t) = e^{-H(t)} = e^{-\int_0^t h(s) ds}$$

$$H(t|x, \beta) = \int_0^t e^{x^T \beta} h_0(u) du = e^{x^T \beta} H_0(t)$$

$$\begin{aligned} \Rightarrow S(t|x, \beta) &= e^{-H(t|x, \beta)} \\ &= e^{-e^{x^T \beta} H_0(t)} = (S_0(t))^{e^{x^T \beta}} \end{aligned}$$

(\Rightarrow not ~~include 1~~ include 1 as the first component)

otherwise we obtain $e^{\beta_0} h_0(t)$, which leads to identifiability problems in this context (we don't know how to estimate separately β_0 and $h_0(t)$)

$$(S_0(t) = e^{-H_0(t)}) \quad \begin{array}{l} \text{survival function} \\ \text{corresponding to} \\ \text{the baseline hazard } h_0 \end{array}$$

Remark:

1. $S_0(t) := e^{-t^\alpha}$ (Weibull($\alpha, 1$)) $\Rightarrow S(t|x, \beta)$ Weibull($\alpha, e^{x^T \beta}$)

if $\check{x}_1^T \beta > \check{x}_2^T \beta \Rightarrow S(t|x_1, \beta) < S(t|x_2, \beta)$ (because $S_0(t) < 1$ and $S(t|x, \beta) = (S_0(t))^{e^{x^T \beta}}$) \Rightarrow recover the AFT model when the error is the Gumbel distribution

Individuals with distinct covariate vectors have ~~the~~ survival function which don't cross

2. x_1, x_2 (\neq): $HR = \frac{h(t|x_1, \beta)}{h(t|x_2, \beta)} = \frac{e^{x_1^T \beta} h_0(t)}{e^{x_2^T \beta} h_0(t)} = e^{(x_1 - x_2)^T \beta} = \text{constant w.r.t. time}$

$$h(t|x_1, \beta) = HR \cdot h(t|x_2, \beta) \quad \text{that's why the name proportional hazard!}$$

Data:

$$D = \{(y_i, \delta_i, x_i)\}_{i=1,\dots,n}$$

$$y_i = \min(T_i, C_i)$$

$$\delta_i = \begin{cases} 1 & \text{if } T_i \leq C_i \\ 0 & \text{if } T_i > C_i \end{cases}$$

Likelihood:

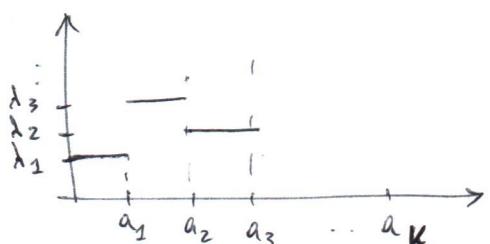
$$L(\beta, h_0 | D) = \prod_{i=1}^n \left(e^{x_i^T \beta} h_0(y_i) \right)^{\delta_i} (S_0(y_i))^{1 - \delta_i} \Rightarrow \text{parameters: } \beta, h_0(\cdot)$$

Priors: $\beta \sim N_k(\underline{0}, \begin{bmatrix} \sigma_1^2 & & \\ & \ddots & \\ & & \sigma_k^2 \end{bmatrix})$ \Rightarrow we assume the parameters to be $\perp\!\!\!\perp$:

$$\beta \sim N_k(\underline{0}, \begin{bmatrix} \sigma_1^2 & & \\ & \ddots & \\ & & \sigma_k^2 \end{bmatrix}) \quad \text{possible assumptions:} \quad \Rightarrow \text{not only } \beta \perp\!\!\!\perp h_0(\cdot), \text{ also } \beta_i \perp\!\!\!\perp \beta_j$$

$$\dim(\beta) = K, \quad \sigma_j \in \text{large value prior on } (\sigma_1^2, \dots, \sigma_K^2)$$

Prior for $h_0(\cdot)$: $\sum_{k=1}^{K+1} \lambda_k \mathbb{1}_{(a_k, a_{k+1})}(t)$ we discretize the function (we split $(0, +\infty)$ into a finite number of subintervals and in each interval $h_0(t)$ is a stepfunction)



We use a finite number of parameters to represent $h_0(t)$:

$$h_0(t) \leftarrow \lambda_1, \lambda_2, \dots, \lambda_{K+1} \rightsquigarrow \text{TF}(\lambda_1, \dots, \lambda_{K+1})$$

~~exchangeable~~

we assume a prior on these λ_i : we want them to be exchangeable, which means:

$$\lambda_i | \gamma \sim \text{gamma}(\gamma, \gamma)$$

$$\gamma \sim \text{gamma}(\alpha_\gamma, \beta_\gamma)$$

Modification of the intrinsic CAR model

Redefine $\Sigma_y^{-1} = D_w - \rho W$ and choose ρ to make Σ_y^{-1} non-singular

- such values for ρ do exist
- replace W by $\tilde{W} = \text{diag}(1/w_{ij})W$; then $\Sigma_y^{-1} = M^{-1}(I - \alpha\tilde{W})$
(M is diagonal) and, if $|\alpha| < 1$, then $I - \alpha\tilde{W}$ is nonsingular

Under $\Sigma_y^{-1} = D_w - \rho W$ (+ symmetry of W):

$$Y_i | y_{-i} \sim \mathcal{N}\left(\rho \frac{\sum_{j=1}^n w_{ij} y_j}{w_{ii}}, \frac{\tau^2}{w_{ii}}\right)$$

Typically $\rho \in (0, 1)$:

$$\rho = 0: Y_i \stackrel{\text{i.i.d.}}{\sim} \mathcal{N}\left(0, \frac{\tau^2}{w_{ii}}\right), i = 1, \dots, n$$

$\rho = 1$: impropriety of $p(y_1, \dots, y_n)$ in (2) (intrinsic CAR model)

prior on ρ in $(0, 1)$: a prior on ρ encouraging spatial association puts mass near 1



GLMM + CAR prior on the spatial random effects

generalized mixed linear models

CARBayes

Study region S partitioned into n non-overlapping areal units

$$S = \{S_1, \dots, S_n\}$$

vector of responses $\mathbf{Y} = (Y_1, \dots, Y_n)$

vector of known offsets $\mathbf{O} = (O_1, \dots, O_n)$

Spatial pattern in the response modelled by

- a matrix of covariates $X = (\mathbf{x}_1^T, \dots, \mathbf{x}_n^T)^T$, where each $\mathbf{x}_k^T = (1, x_{k1}, \dots, x_{kp})$
- a set of random effects $\Phi = (\Phi_1, \dots, \Phi_n)$



GLMM for the likelihood

$$\begin{cases} Y_k | \mu_k \stackrel{\text{ind.}}{\sim} f(y_k | \mu_k, v^2), k = 1, \dots, n \\ g(\mu_k) = \mathbf{x}_k^T \boldsymbol{\beta} + \Phi_k + O_k \end{cases} \quad (3)$$

we model the response $Y_k | \mu_k$ as distributed from the exponential family

$f(y_k | \mu_k, v^2)$ exponential family [Gaussian, binomial, Poisson]

$$E(Y_k) = \mu_k$$

v^2 : scale parameter (when needed)

g : link function [identity, logit, log]

Parameters: $\boldsymbol{\beta}$, v^2 , Φ

$$\boldsymbol{\beta} = (\beta_0, \beta_1, \dots, \beta_p), \Phi = (\Phi_1, \dots, \Phi_n)$$

we're going to model the spatial dependence between the Y_k 's assuming a prior for Φ



which translates somehow the spatial dependence

Models for the likelihood

different likelihoods that the package is able to consider

- Binomial:** $Y_k \sim \text{Binomial}(n_k, \theta_k)$ and
 $\mu_k = \ln(\theta_k / (1 - \theta_k)) = \mathbf{x}_k^T \boldsymbol{\beta} + \Phi_k + O_k$
- Gaussian:** $Y_k \sim N(\mu_k, v^2)$ and $\mu_k = \mathbf{x}_k^T \boldsymbol{\beta} + \Phi_k + O_k$
- Poisson:** $Y_k \sim \text{Poisson}(\mu_k)$ and $\ln(\mu_k) = \mathbf{x}_k^T \boldsymbol{\beta} + \Phi_k + O_k$
- ZIP:** $Y_k \sim \text{ZIP}(\mu_k, \omega_k)$, zero-inflated Poisson model, used to represent data containing an excess of zeros. This is a mixture of a point mass at zero and a Poisson distribution with mean μ_k

O_k is the k th offset: they are typically used when data are counts and we want to model rates (instead of counts), and the time (or area) units are different. For instance, suppose our model is Poisson with $E(Y|x) = \mu_x$ but we assume

$$\log \frac{\mu_x}{t_x} = \beta_0 + \beta_1 x \Rightarrow \log \mu_x = \log(t_x) + \beta_0 + \beta_1 x$$

where t_x is the exposure time for covariate x . Then $\log(t_x)$ is the offset



A priori β , v^2 , Φ are independent

$$\beta \sim \mathcal{N}_{p+1}(0, 1000 I_{p+1}), \quad v^2 \sim \text{inv-gamma}\left(\frac{v_0}{2}, \frac{v_0 \sigma_0^2}{2}\right)$$

$\Phi \sim$ different priors



Independent prior for the random effects

$$\Phi_k | \sigma^2 \stackrel{\text{i.i.d.}}{\sim} \mathcal{N}(0, \sigma^2) \quad \sigma^2 \sim \text{inv-gamma}(a, b) \quad (4)$$

The prior for Φ is exchangeable, since the random effects are conditionally iid (but there is no spatial association included in this prior)

Appropriate if the covariates included in the model have removed all the spatial structure in the response



Global smoothing CAR priors for the random effects

these priors for Φ translates the spatial association

Priors in the CARBayes package:
 (a) intrinsic, (b) BYM models, (c) Leroux et al. (2000) and
 (d) localised spatial autocorrelation by Lee and Mitchell (2012)

we'll focalize on (c)

Each prior is a special case of Gaussian Markov Random Field and can be written as

$$\Phi \sim \mathcal{N}_n(0, \tau^2 \Sigma^{-1}), \text{ i.e.}$$

$$p(\Phi) \propto \exp\left\{-\frac{1}{2\tau^2} \Phi^T (D_w - W) \Phi\right\},$$

with $\Sigma^{-1} = D_w - W$, and W is the proximity matrix.



Global smoothing intrinsic CAR prior (a)

$$\Phi_k | \Phi_{-k}, W, \tau^2 \sim \mathcal{N}\left(\frac{\sum_{j=1}^n w_{kj} \Phi_j}{w_{kk}}, \frac{\tau^2}{w_{kk}}\right) \quad (5)$$

$$\tau^2 \sim \text{inv-gamma}(a = 0.001, b = 0.001)$$

The conditional expectation of each random effect: $E[\Phi_k | \Phi_{-k}, W, \tau^2] =$ average of the random effects in neighboring areas

$\text{Var}[\Phi_k | \Phi_{-k}, W, \tau^2] =$ inversely proportional to the n. of neighbors

Drawbacks:

- (i) the joint prior is improper (but the posterior is proper);
- (ii) it represents strong spatial autocorrelation, and produce random effects that are overly smooth



Most used in practice

Replace Φ_k in (3) with $\Phi_k + \theta_k$, where $(\Phi_1, \dots, \Phi_n) \sim$ as in (5) and $(\theta_1, \dots, \theta_n) \sim$ as in (4)
strong + weak spatial correlation

Identifiability issues: CARBayes estimates their sum $r\theta_k = \Phi_k + \theta_k$



→ Global smoothing intrinsic CAR prior (c)

Alternative CAR prior (Leroux et al., 2000)

$$\begin{aligned} \Phi_k | \Phi_{-k}, W, \tau^2, \rho &\sim \mathcal{N} \left(\frac{\rho \sum_{j=1}^n w_{kj} \Phi_j}{\rho w_{k+} + 1 - \rho}, \frac{\tau^2}{\rho w_{k+} + 1 - \rho} \right) \\ \tau^2 &\sim \text{inv-gamma}(a, b) \\ \rho &\sim \mathcal{U}(0, 1) \end{aligned} \quad (6)$$

Remark

$\rho = 0 \Leftrightarrow$ independence prior
 $\rho = 1 \Leftrightarrow$ intrinsic CAR model prior

This model seems the most appealing from both theoretical and practical standpoints



Global smoothing intrinsic CAR prior (c)

$$\text{Cor}(\Phi_k, \Phi_j | \Phi_{-k,j}, W\rho) = \frac{\rho w_{kj}}{\sqrt{(\rho w_{k+} + 1 - \rho)(\rho w_{j+} + 1 - \rho)}}$$

If $w_{kj} = 0$, i.e. when k and j aren't neighbors, then Φ_k and Φ_j are conditionally independent
If $w_{kj} = 1$, then their partial autocorrelation is controlled by ρ

This prior is overly simplistic since we are assuming a single global level of spatial smoothing for the set of random effects



Localized smoothing CAR priors (d)

They should be able to capture localized spatial autocorrelation, including the identification of boundaries in the random effects surface

Idea: model the elements of W corresponding to geographic adjacent areal units as binary random quantities. Conversely, if (S_k, S_j) do not share a common border, $w_{kj} = 0$

If w_{kj} is estimated as 1 $\Rightarrow \Phi_k$ and Φ_j are spatially correlated, and are smoothed over in the modelling process

If w_{kj} is estimated as 0, then no smoothing is imparted between Φ_k and Φ_j , as they are modelled as conditionally independent. In this case a boundary is said to exist in the random effects surface between areal units (S_k, S_j)



Goal: the aim is to identify the locations of any boundaries (abrupt step changes) in disease risk surfaces, so the available covariates were used to construct dissimilarity metrics rather than being incorporated into the linear predictor.

They model each w_{kj} as a function of the dissimilarity between areal units (S_k, S_j), because large differences in the response are likely to occur where neighboring populations are very different.



CAR priors

Introduce q non-negative dissimilarity metrics $\mathbf{z}_{kj} = (z_{kj1}, \dots, z_{kjq})$, which could include social or physical factors, such as the absolute difference in smoking rates, or the proportion of the shared border that is blocked by a physical barrier (such as a river or railway line) and cannot be crossed.

Binary model:

$$w_{kj}(\boldsymbol{\alpha}) = \begin{cases} 1 & \text{if } \exp(-\sum_{i=1}^q z_{kij}\alpha_i) \geq 0.5 \text{ and } k \sim j \\ 0 & \text{otherwise} \end{cases}$$

$$\alpha_i \stackrel{\text{ind.}}{\sim} \text{Uniform}(0, M_i) \quad \text{for } i = 1, \dots, q$$

The q regression parameters $\boldsymbol{\alpha} = (\alpha_1, \dots, \alpha_q)$ determine the effects of the dissimilarity metrics on $w_{kj}|k \sim j$

For the binary model if $\alpha_i < -\log(0.5)/\max(z_{kij})$, then the i -th dissimilarity metric has not solely identified any boundaries because $\exp(-\alpha_i z_{kij}) > 0.5$ for all $k \sim j$.



CAR priors

$$\mathbf{z}_{kj} = (z_{kj1}, \dots, z_{kjq})$$

Non-binary model:

$$w_{kj}(\boldsymbol{\alpha}) = \exp\left(-\sum_{i=1}^q z_{kij}\alpha_i\right)$$

$$\alpha_i \stackrel{\text{i.i.d.}}{\sim} \text{Uniform}(0, 50) \quad \text{for } i = 1, \dots, q$$

The q regression parameters $\boldsymbol{\alpha} = (\alpha_1, \dots, \alpha_q)$ determine the effects of the dissimilarity metrics on $w_{kj}|k \sim j$



R functions in the CARBayes package

- *S.glm()*: independent random effects
- *S.CARbym()*: BYM model (Gaussian likelihood is not allowed)
- *S.CARleroux()*: CAR in Leroux et al. (2000)
- *S.CARDissimilarity()*: local spatial smoothing
- *S.CARlocalised()*: see the vignette



- ① **Lip cancer data:** how to combine a dataframe and a shapefile together to create an object to run the analysis
 - ② **Property prices in Glasgow:** modelling the spatial pattern in average property prices in Glasgow, to identify the factors that affect property prices and quantify their effects by global smoothing CAR priors
 - ③ **Identifying high-risk disease clusters:** identify boundaries in the risk surface of respiratory disease in Greater Glasgow using a localized smoothing model



Example 3: Identifying high-risk disease clusters

Data: $\mathbf{Y} = (Y_1, \dots, Y_K) \rightarrow$ number of admissions in the hospital due to respiratory diseases

respiratory diseases
E = (E_1, \dots, E_K) → expected value
likelihood

$$Y_k \sim Poisson(E_k B_k), \quad k = 1, \dots, K$$

with R_k is the risk at location S_k

$\log(B_k) = \beta_0 + \phi_k$, $k = 1, \dots, K$ log linear model on the risk

and



Spatio-temporal models

- CARBayesST is a dedicated R package for spatio-temporal areal unit modelling with conditional autoregressive priors
 - Data on a set of K areal units for N consecutive time periods, yielding a rectangular array of $K \times N$ spatio-temporal observations
 - Observations from geographically close areal units and temporally close time periods tend to have more similar values than units and time periods that are further apart
 - the spatio-temporal structure is modelled via set of autocorrelated random effect



Spatio-temporal models

Fit a generalised linear mixed model to these data, whose general form is

$$\begin{aligned} Y_{kt} | \mu_{kt} &\sim f(y_{kt} | \mu_{kt}, v^2), \text{ for } k = 1, 2, \dots, K, t = 1, \dots, N \\ \mu_{kt} &= \mathbf{x}_{kt}^T \boldsymbol{\beta} + O_{kt} + \psi_{kt} \\ \boldsymbol{\beta} &\sim \mathcal{N}(\boldsymbol{\mu}_0, \Sigma_0) \end{aligned}$$

- O_{kt} is an offset: this value is added to the linear predictor of the target (useful in Poisson regression models, where each case may have different levels of exposure to the event of interest)
 - term ψ_{kt} is a latent component for areal unit k and time period t encompassing one or more sets of spatio-temporally autocorrelated random effects: different spatio-temporal structures are available in the package

