

Introduction à la modélisation statistique bayésienne

Un cours en R et Stan avec brms

Ladislav Nalborczyk (LPC, LNC, CNRS, Aix-Marseille Univ)

Planning

Cours n°01 : Introduction à l'inférence bayésienne

Cours n°02 : Modèle Beta-Binomial

Cours n°03 : Introduction à brms, modèle de régression linéaire

Cours n°04 : Modèle de régression linéaire (suite)

Cours n°05 : Markov Chain Monte Carlo

Cours n°06 : Modèle linéaire généralisé

Cours n°07 : Comparaison de modèles

Cours n°08 : Modèles multi-niveaux

Cours n°09 : Modèles multi-niveaux généralisés

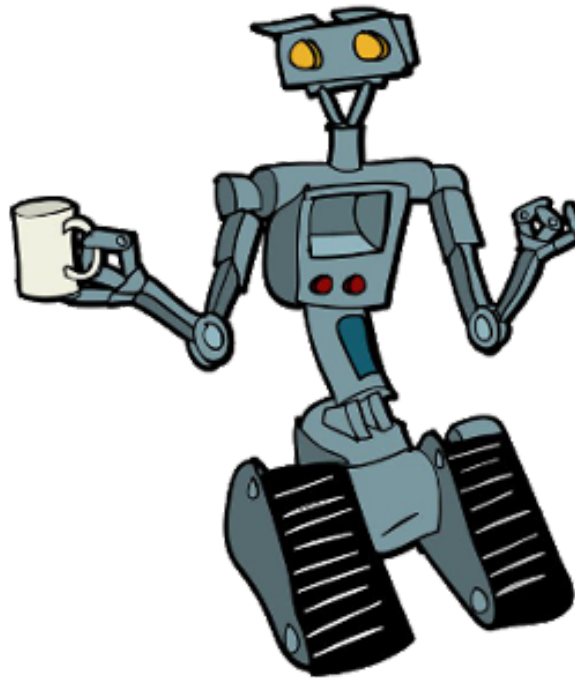
Cours n°10 : Data Hackathon



Modèles multi-niveaux

Le but est de construire un modèle qui puisse **apprendre à plusieurs niveaux**, qui puisse produire des estimations qui seront informées par les différents groupes présents dans les données. Nous allons suivre l'exemple suivant tout au long de ce cours.

Imaginons que nous ayons construit un robot visiteur de cafés, et que celui-ci s'amuse à mesurer le temps d'attente après avoir commandé. Ce robot visite 20 cafés différents, 5 fois le matin et 5 fois l'après-midi, et mesure le temps (en minutes) de service d'un café.



Robot et café

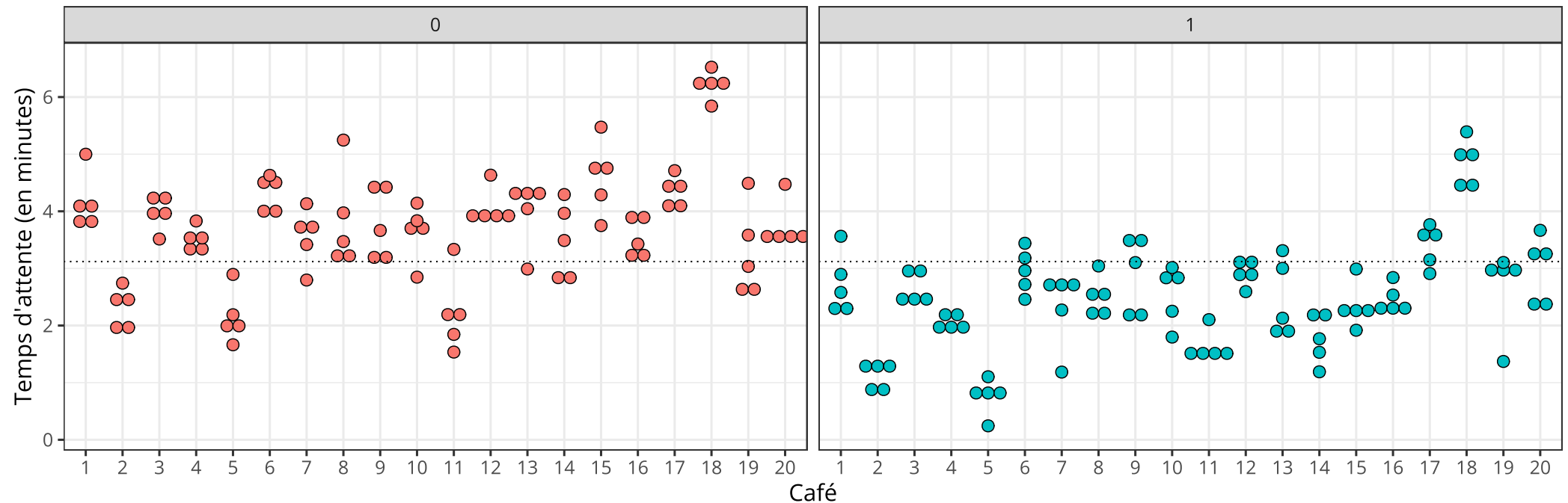
```
1 library(tidyverse)
2 library(imsb)
3
4 df <- open_data(robot)
5 head(x = df, n = 15)
```

	cafe	afternoon	wait
1	1	0	4.9989926
2	1	1	2.2133944
3	1	0	4.1866730
4	1	1	3.5624399
5	1	0	3.9956779
6	1	1	2.8957176
7	1	0	3.7804582
8	1	1	2.3844837
9	1	0	3.8617982
10	1	1	2.5800004
11	2	0	2.7421223
12	2	1	1.3525907
13	2	0	2.5215095
14	2	1	0.9628102
15	2	0	1.9543977



Robot et café

```
1 df %>%
2   ggplot(aes(x = factor(cafe), y = wait, fill = factor(afternoon) )) +
3   geom_dotplot(
4     stackdir = "center", binaxis = "y",
5     dotsize = 1, show.legend = FALSE
6   ) +
7   geom_hline(yintercept = mean(df$wait), linetype = 3) +
8   facet_wrap(~afternoon, ncol = 2) +
9   labs(x = "Café", y = "Temps d'attente (en minutes)")
```



Robot et café, premier modèle

On peut construire un premier modèle, qui estime le temps moyen (sur tous les bistrots confondus) pour être servi.

$$w_i \sim \text{Normal}(\mu_i, \sigma)$$

$$\mu_i = \alpha$$

$$\alpha \sim \text{Normal}(5, 10)$$

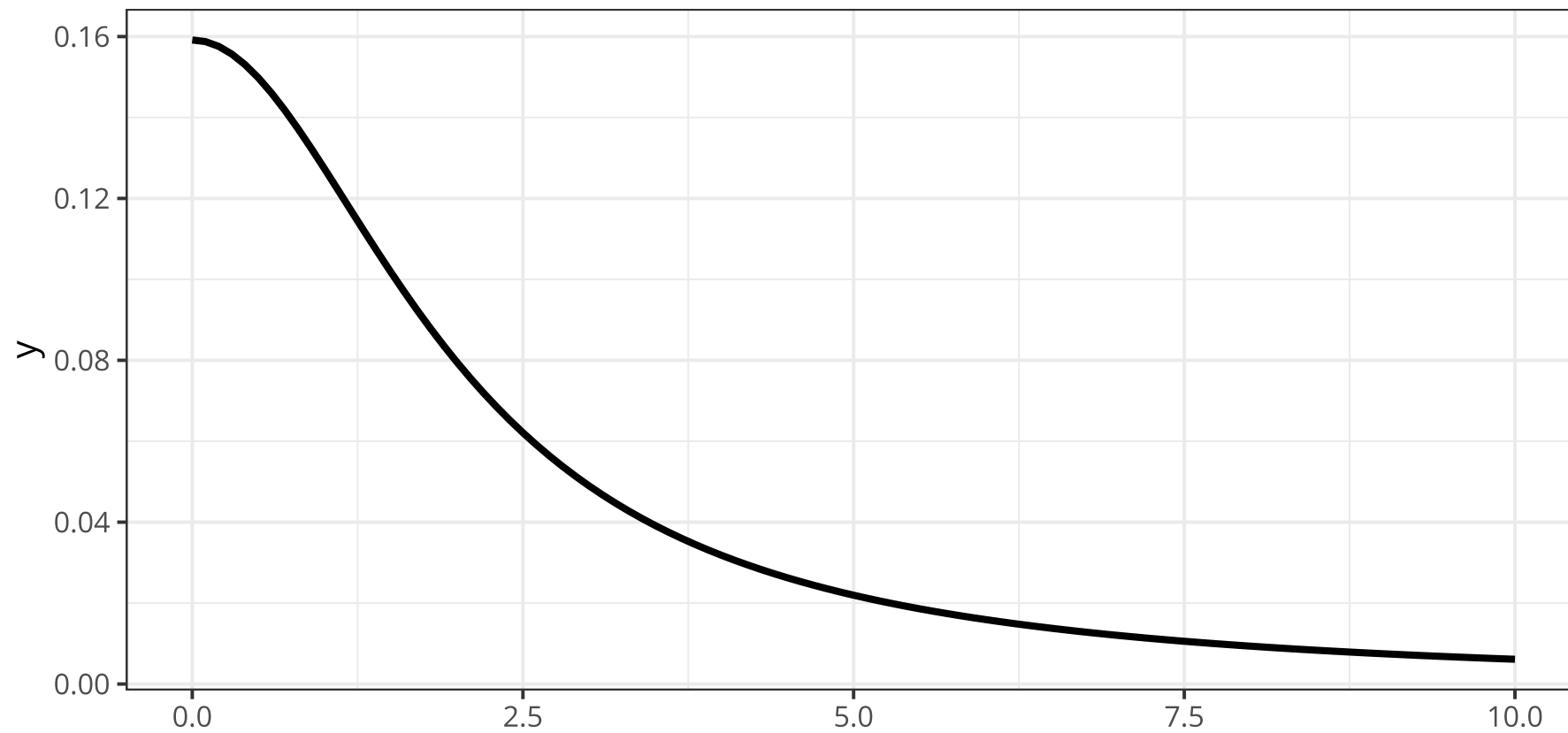
$$\sigma \sim \text{HalfCauchy}(0, 2)$$



Half-Cauchy

$$p(x \mid x_0, \gamma) = \left(\pi \gamma \left[1 + \left(\frac{x - x_0}{\gamma} \right)^2 \right] \right)^{-1}$$

```
1 ggplot(data = data.frame(x = c(0, 10) ), aes(x = x) ) +  
2   stat_function(  
3     fun = dcauchy,  
4     args = list(location = 0, scale = 2), size = 1.5  
5   )
```



Robot et café, premier modèle

```

1 library(brms)
2
3 mod1 <- brm(
4   formula = wait ~ 1,
5   prior = c(
6     prior(normal(5, 10), class = Intercept),
7     prior(cauchy(0, 2), class = sigma)
8   ),
9   data = df,
10  # on utilise tous les coeurs disponibles
11  cores = parallel::detectCores()
12 )

```

```

1 posterior_summary(mod1, probs = c(0.025, 0.975) )

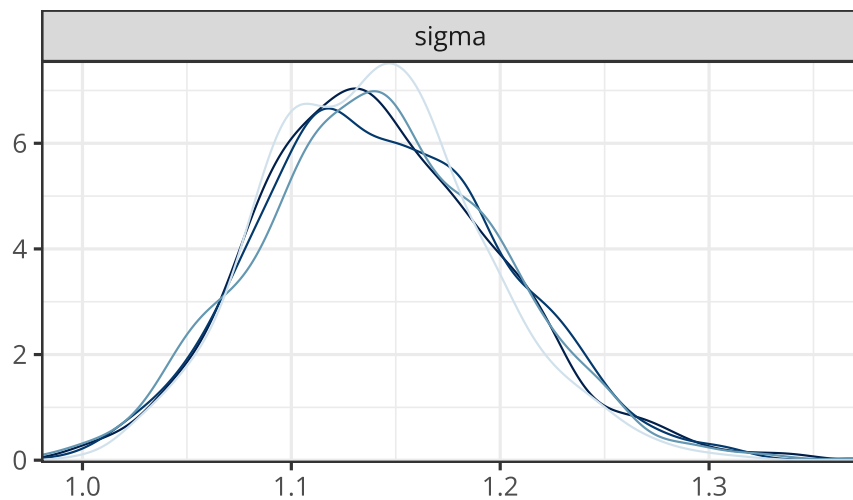
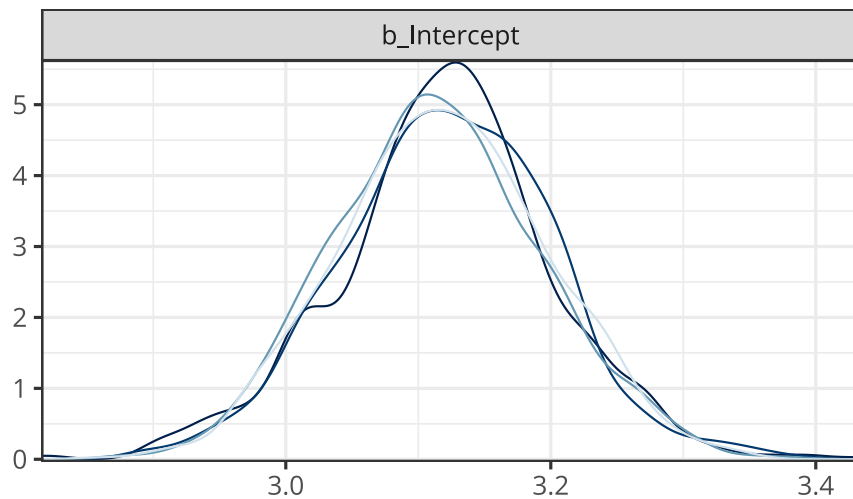
```

	Estimate	Est.Error	Q2.5	Q97.5
b_Intercept	3.118824	0.07953888	2.963235	3.274122
sigma	1.141772	0.05642862	1.038975	1.256931
lprior	-4.666381	0.02438913	-4.716902	-4.622258
lp__	-314.587146	0.98197634	-317.244639	-313.641610



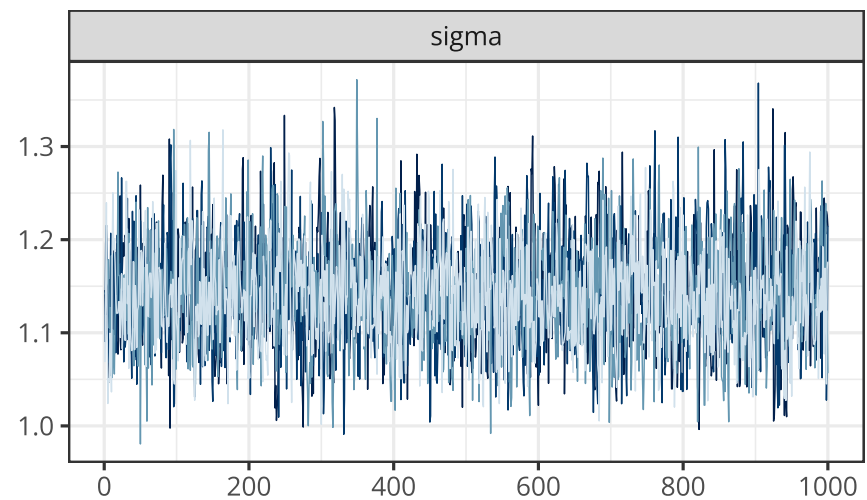
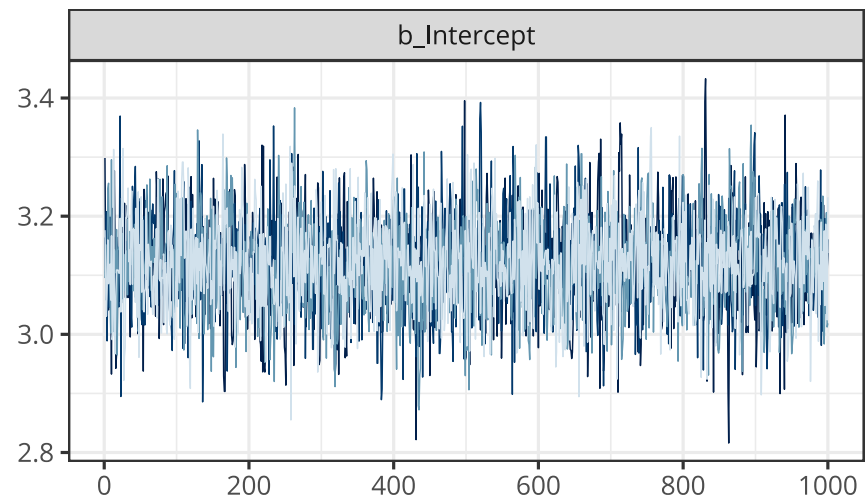
Diagnostic plot

```
1 plot(
2   x = mod1, combo = c("dens_overlay", "trace"),
3   theme = theme_bw(base_size = 16, base_family = "Open Sans")
4 )
```



Chain

- 1
- 2
- 3
- 4



Chain

- 1
- 2
- 3
- 4

Un intercept par café

Deuxième modèle qui estime un intercept par café. Équivalent à construire 20 dummy variables.

$$w_i \sim \text{Normal}(\mu_i, \sigma)$$

$$\mu_i = \alpha_{\text{café}[i]}$$

$$\alpha_{\text{café}[i]} \sim \text{Normal}(5, 10)$$

$$\sigma \sim \text{HalfCauchy}(0, 2)$$

```
1 mod2 <- brm(  
2   formula = wait ~ 0 + factor(cafe),  
3   prior = c(  
4     prior(normal(5, 10), class = b),  
5     prior(cauchy(0, 2), class = sigma)  
6   ),  
7   data = df,  
8   cores = parallel::detectCores()  
9 )
```



Un intercept par café

```
1 posterior_summary(mod2)
```

	Estimate	Est.Error	Q2.5	Q97.5
b_factorcafe1	3.4462051	0.26453027	2.9313403	3.9545029
b_factorcafe2	1.7285618	0.25197656	1.2333122	2.2290702
b_factorcafe3	3.3176327	0.25465844	2.8298730	3.8093179
b_factorcafe4	2.7994639	0.25373088	2.3120310	3.2951321
b_factorcafe5	1.4653843	0.26631968	0.9487625	1.9735277
b_factorcafe6	3.6358637	0.25835749	3.1241016	4.1485547
b_factorcafe7	2.9407022	0.25389813	2.4470846	3.4398958
b_factorcafe8	3.1797369	0.25796189	2.6878408	3.6799060
b_factorcafe9	3.3353910	0.26167094	2.8338354	3.8508551
b_factorcafe10	3.1015820	0.25709832	2.5966855	3.6073649
b_factorcafe11	1.9148883	0.25635381	1.4076090	2.4153084
b_factorcafe12	3.4901065	0.25767956	2.9853369	4.0055029
b_factorcafe13	3.2202182	0.26072015	2.7211301	3.7359100
b_factorcafe14	2.6301540	0.25548824	2.1134215	3.1486889
b_factorcafe15	3.4806255	0.25638191	2.9866802	3.9839562
b_factorcafe16	3.0002179	0.27003798	2.4678319	3.5217938
b_factorcafe17	3.8749534	0.26143587	3.3764925	4.3908258
b_factorcafe18	5.5303247	0.26354702	5.0243866	6.0329810
b_factorcafe19	2.9747747	0.26608316	2.4568977	3.4897951
b_factorcafe20	3.3645093	0.25942162	2.8460383	3.8725111



Modèle multi-niveaux

Est-ce qu'on ne pourrait pas faire en sorte que le temps mesuré au café 1 **informe** la mesure réalisée au café 2, et au café 3 ? Ainsi que le temps moyen pour être servi ? Nous allons apprendre le prior à partir des données...

Niveau 1 : $w_i \sim \text{Normal}(\mu_i, \sigma)$

$$\mu_i = \alpha_{\text{café}[i]}$$

Niveau 2 : $\alpha_{\text{café}} \sim \text{Normal}(\alpha, \sigma_{\text{café}})$

$$\alpha \sim \text{Normal}(5, 10)$$

$$\sigma_{\text{café}} \sim \text{HalfCauchy}(0, 2)$$

$$\sigma \sim \text{HalfCauchy}(0, 2)$$

Le prior de l'intercept pour chaque café ($\alpha_{\text{café}}$) est maintenant fonction de deux paramètres (α et $\sigma_{\text{café}}$). α et $\sigma_{\text{café}}$ sont appelés des **hyper-paramètres**, ce sont des paramètres pour des paramètres, et leurs priors sont appelés des **hyperpriors**. Il y a deux niveaux dans le modèle...



Équivalences (encore)

$$w_i \sim \text{Normal}(\mu_i, \sigma)$$

$$\mu_i = \alpha_{\text{café}[i]}$$

$$\alpha_{\text{café}} \sim \text{Normal}(\alpha, \sigma_{\text{café}})$$

NB : α est ici défini dans le prior de $\alpha_{\text{café}}$ mais on pourrait, de la même manière, le définir dans le modèle linéaire :

$$w_i \sim \text{Normal}(\mu_i, \sigma)$$

$$\mu_i = \alpha + \alpha_{\text{café}[i]}$$

$$\alpha_{\text{café}} \sim \text{Normal}(0, \sigma_{\text{café}})$$

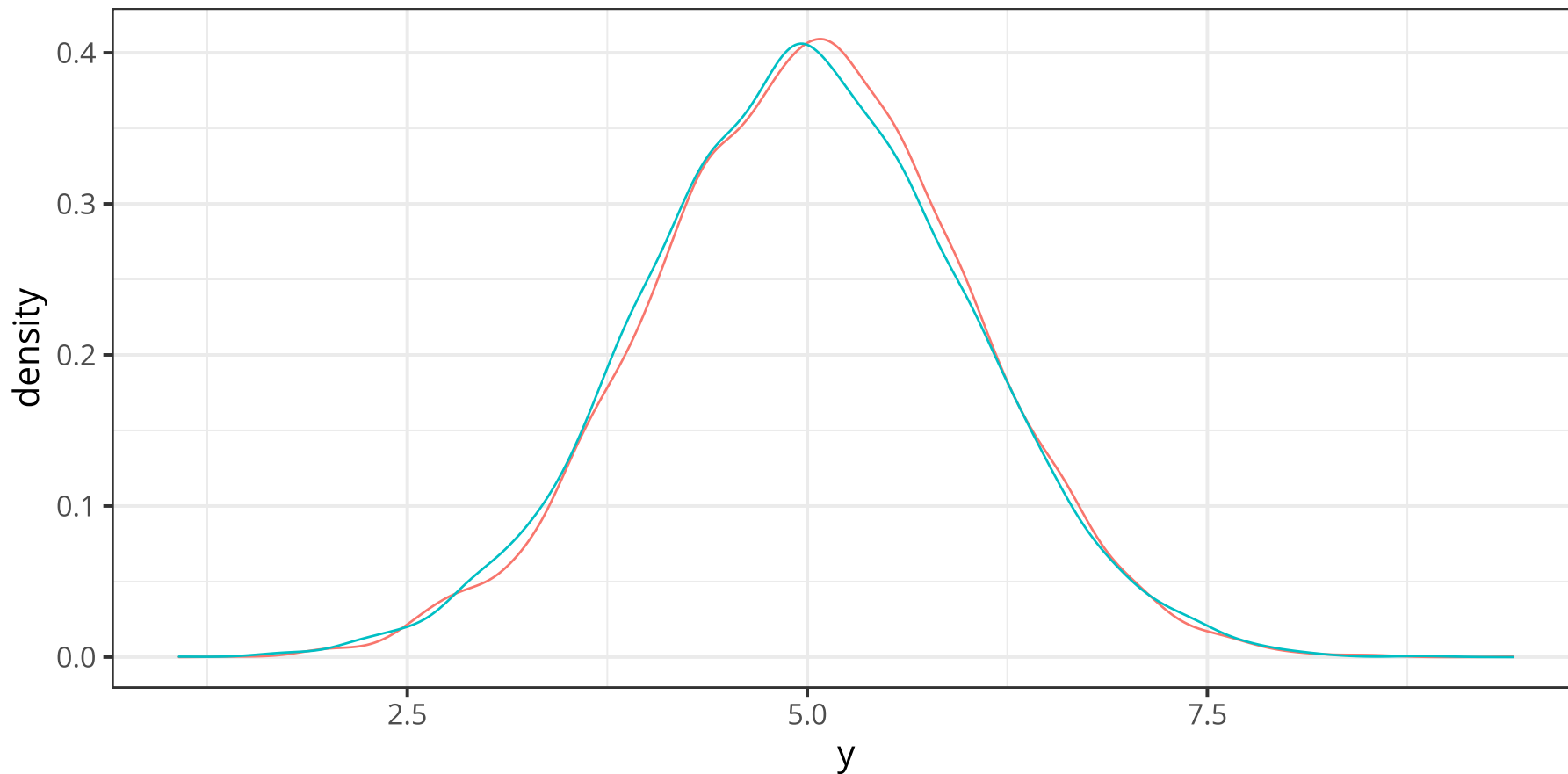
On peut toujours “enlever” la moyenne d’une distribution gaussienne et la considérer comme une constante plus une gaussienne centrée sur zéro.

NB : quand α est défini dans le modèle linéaire, les $\alpha_{\text{café}}$ représentent des déviations de l’intercept moyen. Il faut donc ajouter α et $\alpha_{\text{café}}$ pour obtenir le temps d’attente moyen par café...



Équivalences (encore)

```
1 y1 <- rnorm(n = 1e4, mean = 5, sd = 1)
2 y2 <- rnorm(n = 1e4, mean = 0, sd = 1) + 5
3
4 data.frame(y1 = y1, y2 = y2) %>%
5   pivot_longer(cols = 1:2, names_to = "x", values_to = "y") %>%
6   ggplot(aes(x = y, colour = x)) +
7   geom_density(show.legend = FALSE)
```





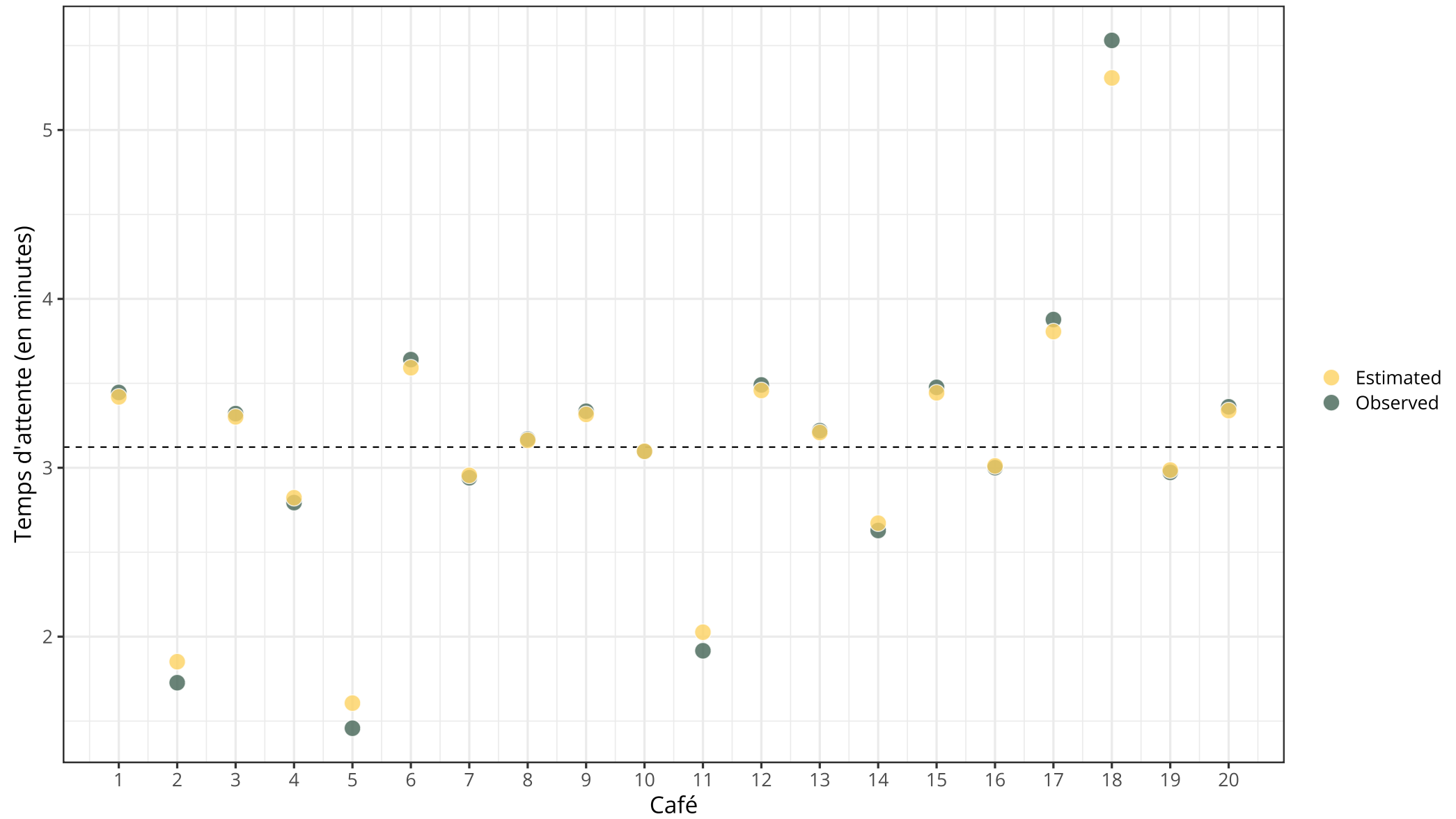
Modèle multi-niveaux

```
1 mod3 <- brm(  
2   formula = wait ~ 1 + (1 | cafe),  
3   prior = c(  
4     prior(normal(5, 10), class = Intercept),  
5     prior(cauchy(0, 2), class = sigma),  
6     prior(cauchy(0, 2), class = sd)  
7   ),  
8   data = df,  
9   warmup = 1000, iter = 5000,  
10  cores = parallel::detectCores()  
11  )
```

Ce modèle a 23 paramètres, l'intercept général α , la variabilité résiduelle σ , la variabilité entre les cafés $\sigma_{\text{café}}$, et un intercept par café.



Shrinkage



Shrinkage magic (Efron & Morris, 1977)

Stein's Paradox in Statistics

The best guess about the future is usually obtained by computing the average of past events. Stein's paradox defines circumstances in which there are estimators better than the arithmetic average

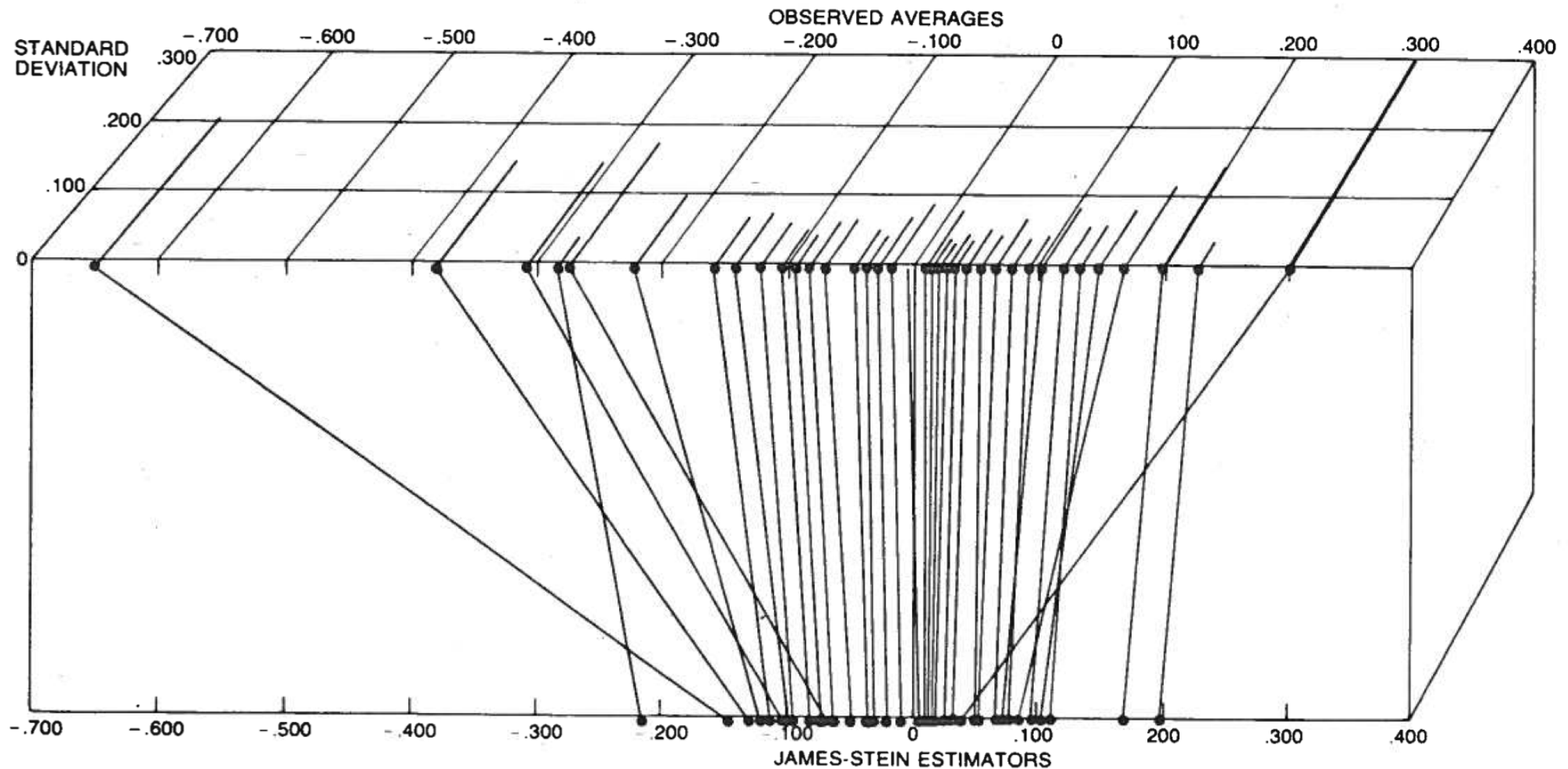
by Bradley Efron and Carl Morris

L'estimateur James-Stein est défini comme $z = \bar{y} + c(y - \bar{y})$, où \bar{y} désigne la moyenne de l'échantillon, y une observation individuelle, et c une constante, le **shrinking factor** (Efron & Morris, 1977).



Shrinkage magic (Efron & Morris, 1977)

Le shrinking factor est déterminé à la fois par la variabilité (imprécision) de la mesure (e.g., son écart-type) et par la distance à l'estimation moyenne (i.e., $y - \bar{y}$). En d'autres termes, cet estimateur fait moins "confiance" (i.e., accorde moins de poids) aux observations imprécises et/ou extrêmes. Le shrinkage agit donc comme une protection contre le sur-apprentissage (overfitting).



Pooling

Le **shrinkage** observé slide précédente est dû à des phénomènes de partage (pooling) de l'information entre les cafés. L'estimation de l'intercept pour chaque café informe l'estimation de l'intercept des autres cafés, ainsi que l'estimation de l'intercept général (i.e., la moyenne générale).

On distingue en général trois perspectives (ou stratégies) :

- **Complete pooling** : on suppose que le temps d'attente est invariant, on estime un intercept commun (mod1).
- **No pooling** : on suppose que les temps d'attente de chaque café sont uniques et indépendants: on estime un intercept par café, mais sans informer le **niveau** supérieur (mod2).
- **Partial pooling** : on utilise un prior adaptatif, comme dans l'exemple précédent (mod3).

La stratégie **complete pooling** en général underfitte les données (faibles capacités de prédiction) tandis que le stratégie **no pooling** revient à overfitter les données (faibles capacités de prédiction ici aussi). La stratégie **partial pooling** (i.e., celle des modèles multi-niveaux) permet d'équilibrer underfitting et overfitting.



Comparaison de modèles

On peut comparer ces trois modèles en utilisant le WAIC (discuté au Cours n°07).

```
1 # calcul du WAIC et ajout du WAIC à chaque modèle
2 mod1 <- add_criterion(mod1, "waic")
3 mod2 <- add_criterion(mod2, "waic")
4 mod3 <- add_criterion(mod3, "waic")
5
6 # comparaison des WAIC de chaque modèle
7 w <- loo_compare(mod1, mod2, mod3, criterion = "waic")
8 print(w, simplify = FALSE)
```

	elpd_diff	se_diff	elpd_waic	se_elpd_waic	p_waic	se_p_waic	waic	se_waic
mod3	0.0	0.0	-253.8	8.3	18.1	1.5	507.6	16.6
mod2	-0.6	1.3	-254.4	8.4	19.4	1.6	508.7	16.7
mod1	-57.3	10.6	-311.1	10.6	2.0	0.3	622.1	21.1

On remarque que le modèle 3 a seulement 18 “effective parameters” (p_{WAIC}), et moins de paramètres que le modèle 2, alors qu’il en a en réalité 2 de plus... `posterior_summary(mod3)[3, 1]` nous donne le sigma du prior adaptatif des $\alpha_{\text{café}}$ ($\sigma_{\text{café}} = 0.82$). On remarque que ce sigma est très faible et correspond à assigner un prior très contraignant, ou **régularisateur**.



Comparaison de modèles

On compare les estimations du premier modèle (complete pooling model) et du troisième modèle (partial pooling model).

```
1 posterior_summary(mod1, pars = c("^b", "sigma") )
```

	Estimate	Est.Error	Q2.5	Q97.5
b_Intercept	3.118824	0.07953888	2.963235	3.274122
sigma	1.141772	0.05642862	1.038975	1.256931

```
1 posterior_summary(mod3, pars = c("^b", "sigma") )
```

	Estimate	Est.Error	Q2.5	Q97.5
b_Intercept	3.1224462	0.20548659	2.7134453	3.5301519
sigma	0.8219151	0.04328802	0.7430558	0.9108058

Les deux modèles font la même prédiction (en moyenne) pour α , mais le modèle 3 est plus incertain de sa prédiction que le modèle 1 (voir l'erreur standard pour α)...

L'estimation de σ du modèle 3 est plus petite que celle du modèle 1 car le modèle 3 **décompose** la variabilité non expliquée en deux sources : la variabilité du temps d'attente entre les cafés $\sigma_{\text{café}}$ et la variabilité résiduelle σ .



Robot et café

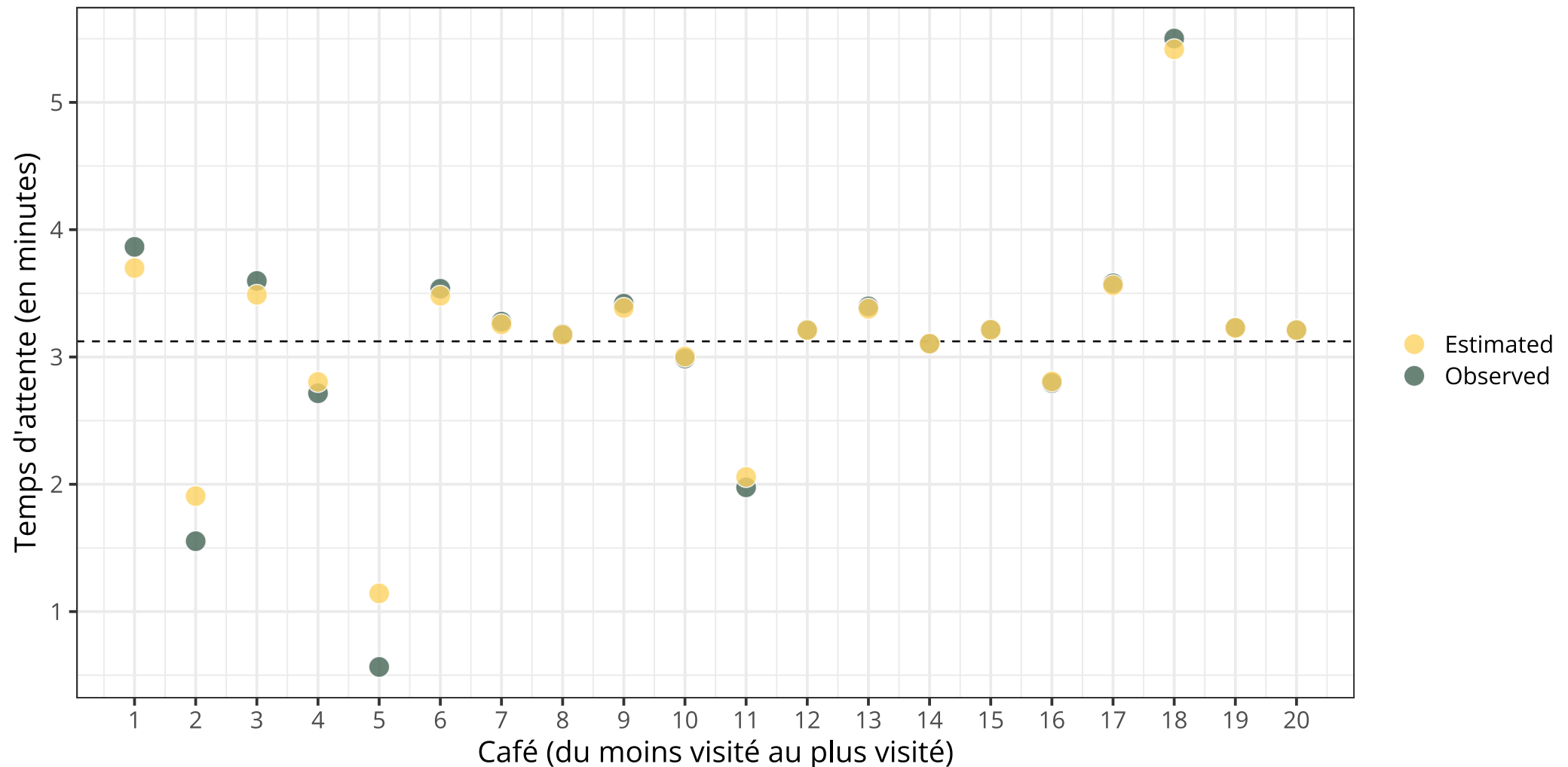
Imaginons que notre robot ne visite pas tous les cafés le même nombre de fois (comme dans le cas précédent) mais qu'il visite plus souvent les cafés proches de chez lui...

```
1 df2 <- open_data(robot_unequal) # nouveau jeu de données
2
3 mod4 <- brm(
4   formula = wait ~ 1 + (1 | cafe),
5   prior = c(
6     prior(normal(5, 10), class = Intercept),
7     prior(cauchy(0, 2), class = sigma),
8     prior(cauchy(0, 2), class = sd)
9   ),
10  data = df2,
11  warmup = 1000, iter = 5000,
12  cores = parallel::detectCores()
13 )
```



Shrinkage

On observe que les cafés qui sont souvent visités (à droite) subissent moins l'effet du **shrinkage**. Leur estimation est moins “tirée” vers la moyenne générale que les estimations des cafés les moins souvent visités (à gauche).



Aparté : effets fixes et effets aléatoires

Cinq définitions (contradictaires) relevées par Gelman ([2005](#)).

- Fixed effects are constant across individuals, and random effects vary.
- Effects are fixed if they are interesting in themselves or random if there is interest in the underlying population.
- When a sample exhausts the population, the corresponding variable is fixed; when the sample is a small (i.e., negligible) part of the population the corresponding variable is random.
- If an effect is assumed to be a realized value of a random variable, it is called a random effect.
- Fixed effects are estimated using least squares (or, more generally, maximum likelihood) and random effects are estimated with shrinkage.

Gelman & Hill ([2006](#)) suggèrent plutôt l'utilisation des termes de **constant effects** et **varying effects**, et de toujours utiliser la modélisation multi-niveaux, en considérant que ce qu'on appelle **effet fixe** peut simplement être considéré comme un **effet aléatoire** dont la variance serait égale à 0.



Régularisation et terminologie

Le fait de faire varier les intercepts de chaque café est simplement une autre manière de régulariser (de manière adaptative), c'est à dire de diminuer le poids accordé aux données dans l'estimation. Le modèle devient à même d'estimer à quel point les groupes (ici les cafés) sont différents, tout en estimant les caractéristiques de chaque café...

Différence entre les **cross-classified** (ou “crossed”) multilevel models et **nested or hierarchical** multilevel models. Le premier type de modèle concerne des données structurées selon deux (ou plus) facteurs aléatoires non “nichés”. Le deuxième type de modèles concerne des données structurées de manière hiérarchique (e.g., un élève dans une classe dans une école dans une ville...). Voir [cette discussion](#) pour plus de détails.

Les deux types de modèles s'écrivent cependant de manière similaire, sur plusieurs “niveaux”. Le terme “multi-niveaux” (dans notre terminologie) fait donc référence à la structure du modèle, à sa spécification. À distinguer de la structure des données.



Exemple de modèle “cross-classified”

On pourrait se poser la question de savoir si la récence des cafés (leur âge) ne serait pas une source de variabilité non contrôlée ? Il suffit d'ajouter un intercept qui varie par âge, et de lui attribuer un prior adaptatif.

$$w_i \sim \text{Normal}(\mu_i, \sigma)$$

$$\mu_i = \alpha + \alpha_{\text{café}[i]} + \alpha_{\text{âge}[i]}$$

$$\alpha_{\text{café}} \sim \text{Normal}(5, \sigma_{\text{café}})$$

$$\alpha_{\text{âge}} \sim \text{Normal}(5, \sigma_{\text{âge}})$$

$$\alpha \sim \text{Normal}(0, 10)$$

$$\sigma_{\text{café}} \sim \text{HalfCauchy}(0, 2)$$

$$\sigma_{\text{âge}} \sim \text{HalfCauchy}(0, 2)$$

$$\sigma \sim \text{HalfCauchy}(0, 2)$$



Robot et café : varying intercept + varying slope

On s'intéresse maintenant à l'effet du moment de la journée sur le temps d'attente. Attend-on plus le matin, ou l'après-midi ?

$$w_i \sim \text{Normal}(\mu_i, \sigma)$$

$$\mu_i = \alpha_{\text{café}[i]} + \beta_{\text{café}[i]} \times A_i$$

Où A_i est une dummy variable codée 0/1 pour le matin et l'après-midi et où $\beta_{\text{café}}$ est donc un paramètre de différence (i.e., une pente) entre le matin et l'après-midi.

Remarque : on sait que les cafés ont des intercepts et des pentes qui co-varient... Les cafés populaires seront surchargés le matin et beaucoup moins l'après-midi, résultant en une pente importante. Ces cafés auront aussi un temps d'attente moyen plus long (i.e., un intercept plus grand). Dans ces cafés, α est grand et β est loin de zéro. À l'inverse, dans un café peu populaire, le temps d'attente sera faible, ainsi que la différence entre matin et après-midi.

On pourrait donc utiliser la co-variation entre intercept et pente pour faire de meilleures inférences. Autrement dit, faire en sorte que l'estimation de l'intercept informe celle de la pente, et réciproquement.



Robot et café : varying intercept + varying slope

On s'intéresse maintenant à l'effet du moment de la journée sur le temps d'attente. Attend-on plus le matin, ou l'après-midi ?

$$\begin{aligned} w_i &\sim \text{Normal}(\mu_i, \sigma) \\ \mu_i &= \alpha_{\text{café}[i]} + \beta_{\text{café}[i]} \times A_i \\ \begin{bmatrix} \alpha_{\text{café}} \\ \beta_{\text{café}} \end{bmatrix} &\sim \text{MVNormal}\left(\begin{bmatrix} \alpha \\ \beta \end{bmatrix}, \mathbf{S}\right) \end{aligned}$$

La troisième ligne postule que chaque café a un intercept $\alpha_{\text{café}}$ et une pente $\beta_{\text{café}}$, définis par un prior Gaussien bivarié (i.e., à deux dimensions) ayant comme moyennes α et β et comme matrice de covariance \mathbf{S} .



Aparté : distribution gaussienne multivariée

$$\mathbf{x} \sim \mathcal{N}(\boldsymbol{\mu}, \boldsymbol{\Sigma})$$

Où $\boldsymbol{\mu}$ est un vecteur (à k dimensions) de moyennes, par exemple: `mu <- c(a, b)`.

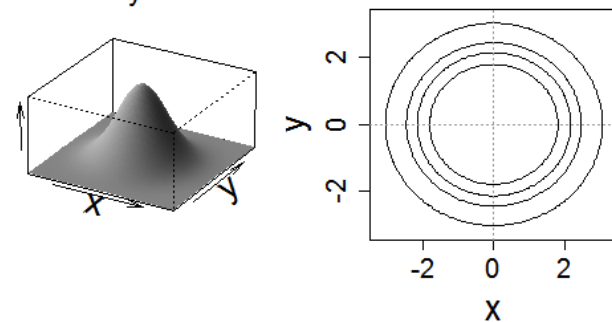
$\boldsymbol{\Sigma}$ est une matrice de covariance de $k \times k$ dimensions, et qui correspond à la matrice donnée par la fonction `vcov()`.

$$\boldsymbol{\Sigma} = \begin{pmatrix} \sigma_{\alpha}^2 & \sigma_{\alpha}\sigma_{\beta}\rho \\ \sigma_{\alpha}\sigma_{\beta}\rho & \sigma_{\beta}^2 \end{pmatrix}$$

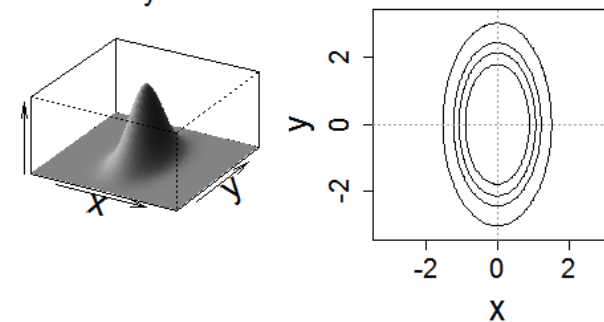


Aparté : distribution gaussienne multivariée

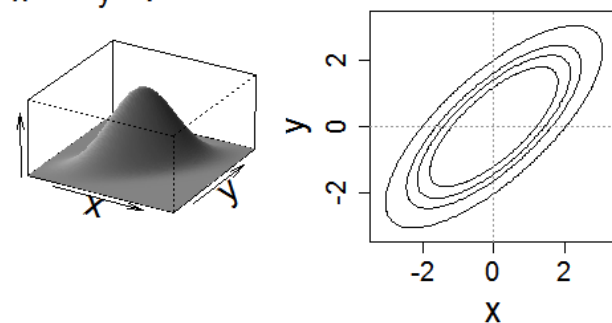
$$\sigma_x = \sigma_y, \rho = 0$$



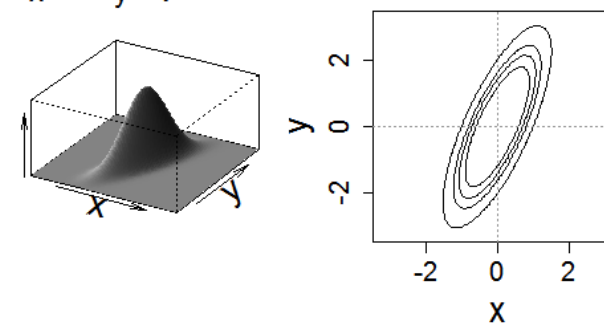
$$2\sigma_x = \sigma_y, \rho = 0$$



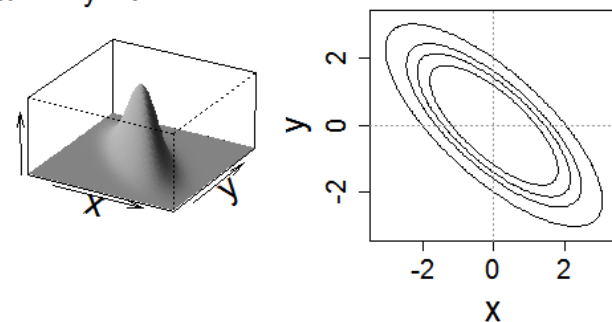
$$\sigma_x = \sigma_y, \rho = 0.75$$



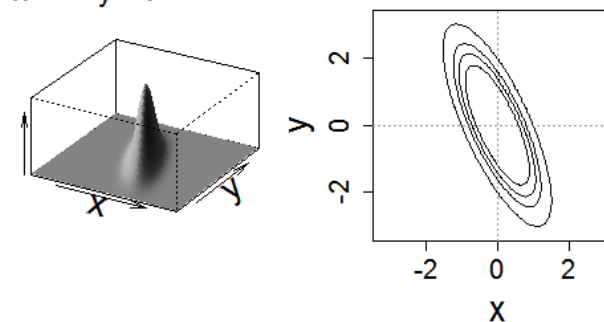
$$2\sigma_x = \sigma_y, \rho = 0.75$$



$$\sigma_x = \sigma_y, \rho = -0.75$$



$$2\sigma_x = \sigma_y, \rho = -0.75$$



Aparté : distribution gaussienne multivariée

$$\Sigma = \begin{pmatrix} \sigma_\alpha^2 & \sigma_\alpha \sigma_\beta \rho \\ \sigma_\alpha \sigma_\beta \rho & \sigma_\beta^2 \end{pmatrix}$$

Cette matrice peut se construire de deux manières différentes, strictement équivalentes.

```
1 sigma_a <- 1
2 sigma_b <- 0.75
3 rho <- 0.7
4 cov_ab <- sigma_a * sigma_b * rho
5 (Sigma1 <- matrix(c(sigma_a^2, cov_ab, cov_ab, sigma_b^2), ncol = 2) )
```

```
      [,1] [,2]
[1,] 1.000 0.5250
[2,] 0.525 0.5625
```



Aparté : distribution gaussienne multivariée

$$\Sigma = \begin{pmatrix} \sigma_{\alpha}^2 & \sigma_{\alpha}\sigma_{\beta}\rho \\ \sigma_{\alpha}\sigma_{\beta}\rho & \sigma_{\beta}^2 \end{pmatrix}$$

La deuxième méthode est pratique car elle considère séparément les écart-types et les corrélations.

```
1 (sigmas <- c(sigma_a, sigma_b) ) # standard deviations
```

```
[1] 1.00 0.75
```

```
1 (Rho <- matrix(c(1, rho, rho, 1), nrow = 2) ) # correlation matrix
```

```
      [,1] [,2]
[1,]  1.0  0.7
[2,]  0.7  1.0
```

```
1 (Sigma2 <- diag(sigmas) %*% Rho %*% diag(sigmas) )
```

```
      [,1] [,2]
[1,] 1.000 0.5250
[2,] 0.525 0.5625
```



Robot et café : varying intercept + varying slope

$$w_i \sim \text{Normal}(\mu_i, \sigma)$$

$$\mu_i = \alpha_{\text{café}[i]} + \beta_{\text{café}[i]} \times A_i$$

$$\begin{bmatrix} \alpha_{\text{café}} \\ \beta_{\text{café}} \end{bmatrix} \sim \text{MVNormal} \left(\begin{bmatrix} \alpha \\ \beta \end{bmatrix}, \mathbf{S} \right)$$

$$\mathbf{S} = \begin{pmatrix} \sigma_\alpha & 0 \\ 0 & \sigma_\beta \end{pmatrix} \mathbf{R} \begin{pmatrix} \sigma_\alpha & 0 \\ 0 & \sigma_\beta \end{pmatrix}$$

$$\alpha \sim \text{Normal}(0, 10)$$

$$\beta \sim \text{Normal}(0, 10)$$

$$\sigma_\alpha \sim \text{HalfCauchy}(0, 2)$$

$$\sigma_\beta \sim \text{HalfCauchy}(0, 2)$$

$$\sigma \sim \text{HalfCauchy}(0, 2)$$

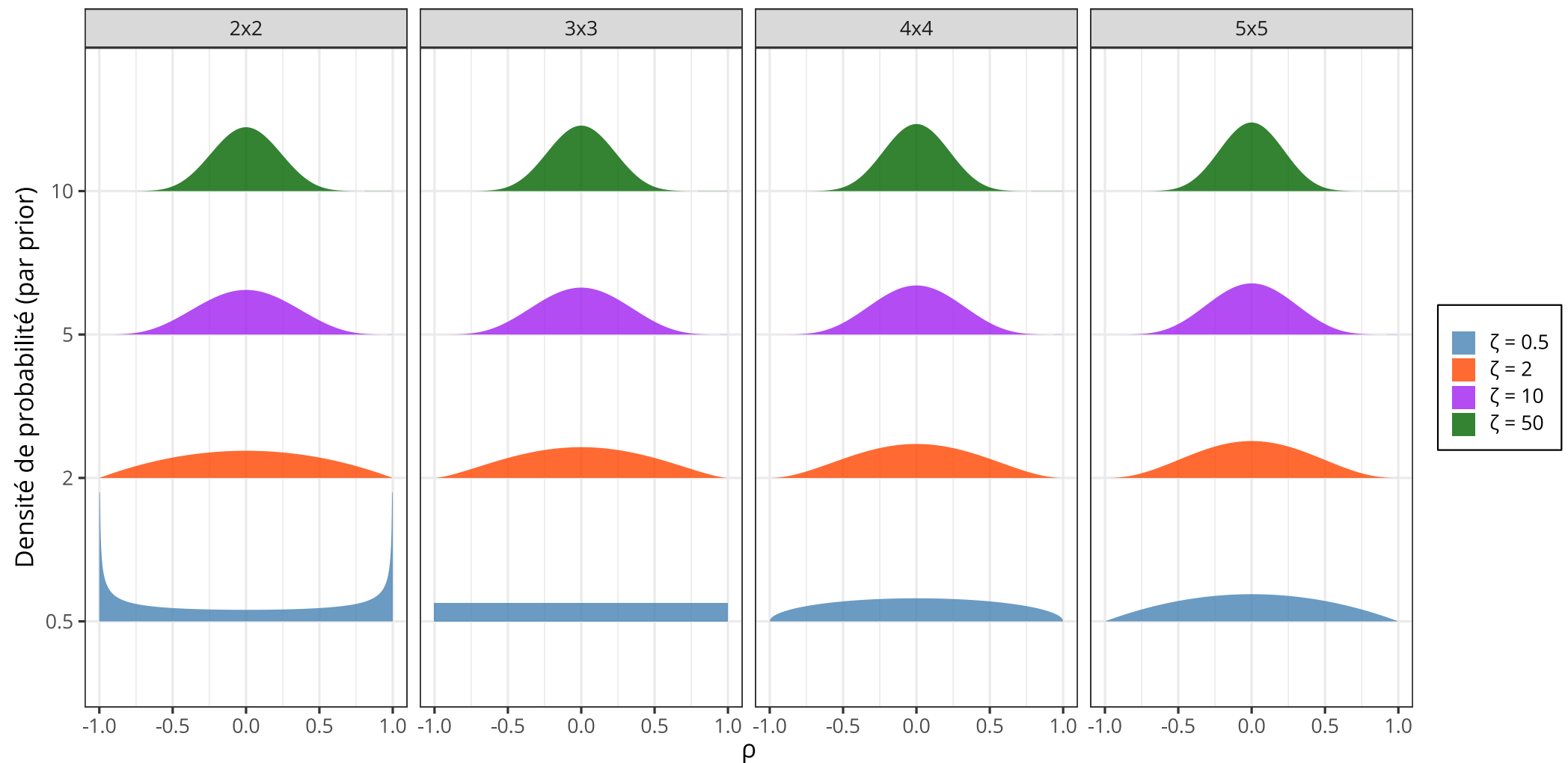
$$\mathbf{R} \sim \text{LKJ}(2)$$

\mathbf{S} est définie en factorisant σ_α , σ_β , et la matrice de corrélation \mathbf{R} . La suite du modèle définit simplement les priors pour les effets constants. La dernière ligne spécifie le prior pour \mathbf{R} .



LKJ prior

Prior proposé par Lewandowski et al. ([2009](#)). Un seul paramètre ζ (zeta) spécifie la concentration de la distribution du coefficient de corrélation. Le prior **LKJ(2)** définit un prior peu informatif pour ρ (rho) qui est sceptique des corrélations extrêmes (i.e., des valeurs proches de -1 ou 1).



Rappels de syntaxe

Le package `brms` utilise la même syntaxe que les fonctions de base R (comme `lm`) ou que le package `lme4`.

```
1 Reaction ~ Days + (1 + Days | Subject)
```

La partie gauche représente notre variable dépendante (ou “outcome”, i.e., ce qu’on essaye de prédire).

La partie droite permet de définir les prédicteurs. L’intercept est généralement implicite, de sorte que les deux écritures ci-dessous sont équivalentes.

```
1 c(Reaction, Memory) ~ Days + (1 + Days | Subject)
2 c(Reaction, Memory) ~ 1 + Days + (1 + Days | Subject)
```



Rappels de syntaxe

La première partie de la partie droite de la formule représente les effets constants (effets fixes), tandis que la seconde partie (entre parenthèses) représente les effets “variants” ou “variables” (effets aléatoires).

```
1 c(Reaction, Memory) ~ Days + (1 | Subject)
2 c(Reaction, Memory) ~ Days + (Days | Subject)
```



Le premier modèle ci-dessus contient seulement un intercept variable, qui varie par **Subject**. Le deuxième modèle contient également un intercept variable, mais aussi une pente variable pour l'effet de **Days**.



Rappels de syntaxe

Lorsqu'on inclut plusieurs effets variants (e.g., un intercept et une pente variables), **brms** postule qu'on souhaite aussi estimer la corrélation entre ces deux effets. Dans le cas contraire, on peut supprimer cette corrélation (i.e., la fixer à 0) en utilisant `||`.

```
1 c(Reaction, Memory) ~ Days + (1 + Days || Subject)
```

Les modèles précédents postulaient un modèle génératif Gaussien. Ce postulat peut être changé facilement en spécifiant la fonction souhaitée via l'argument **family**.

```
1 brm(Reaction ~ 1 + Days + (1 + Days | Subject), family = lognormal() )
```



Modèle brms

On spécifie un intercept et une pente (pour l'effet d'**afternoon**) qui varient par **cafe**.

```
1 mod5 <- brm(  
2   wait ~ 1 + afternoon + (1 + afternoon | cafe),  
3   prior = c(  
4     prior(normal(0, 10), class = Intercept),  
5     prior(normal(0, 10), class = b),  
6     prior(cauchy(0, 2), class = sigma),  
7     prior(cauchy(0, 2), class = sd)  
8   ),  
9   data = df,  
10  warmup = 1000, iter = 5000,  
11  cores = parallel::detectCores()  
12 )
```

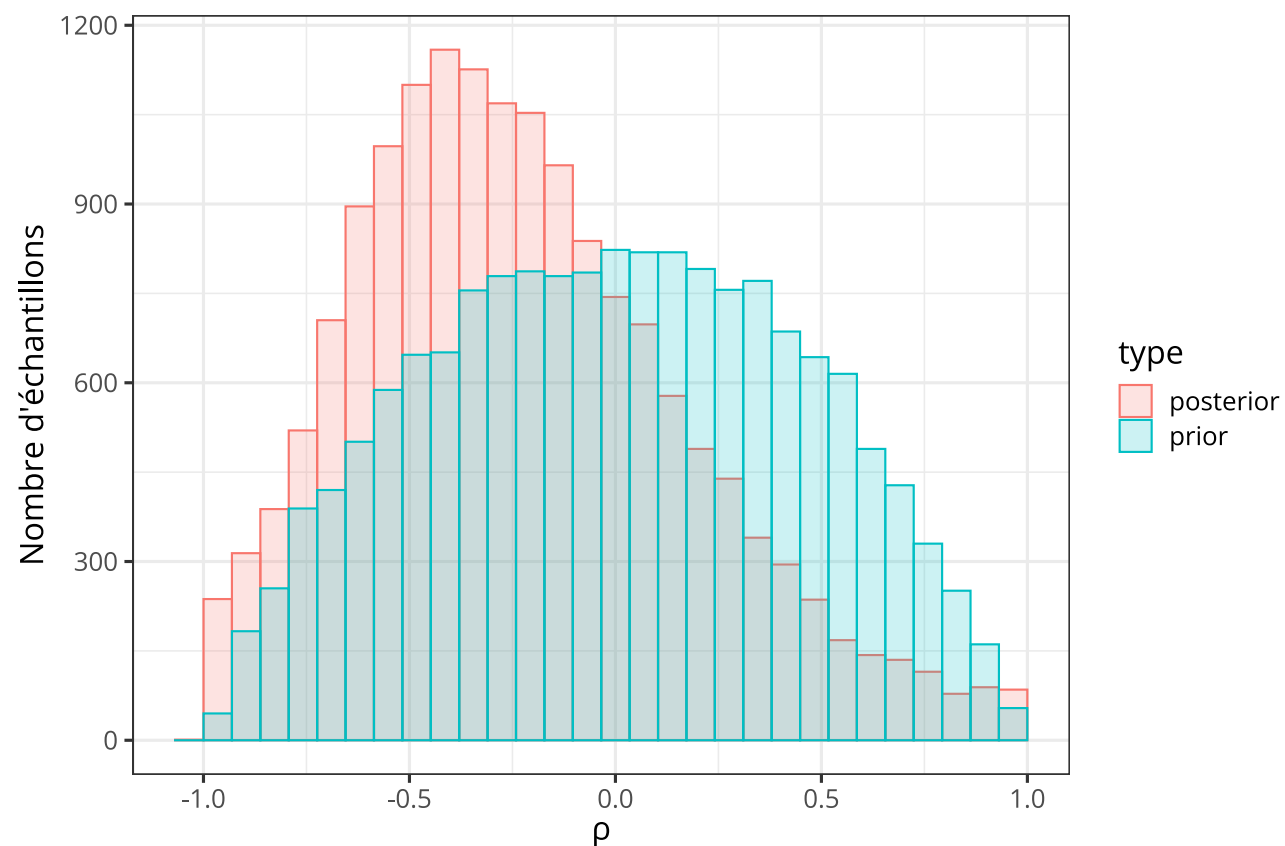


Distribution postérieure

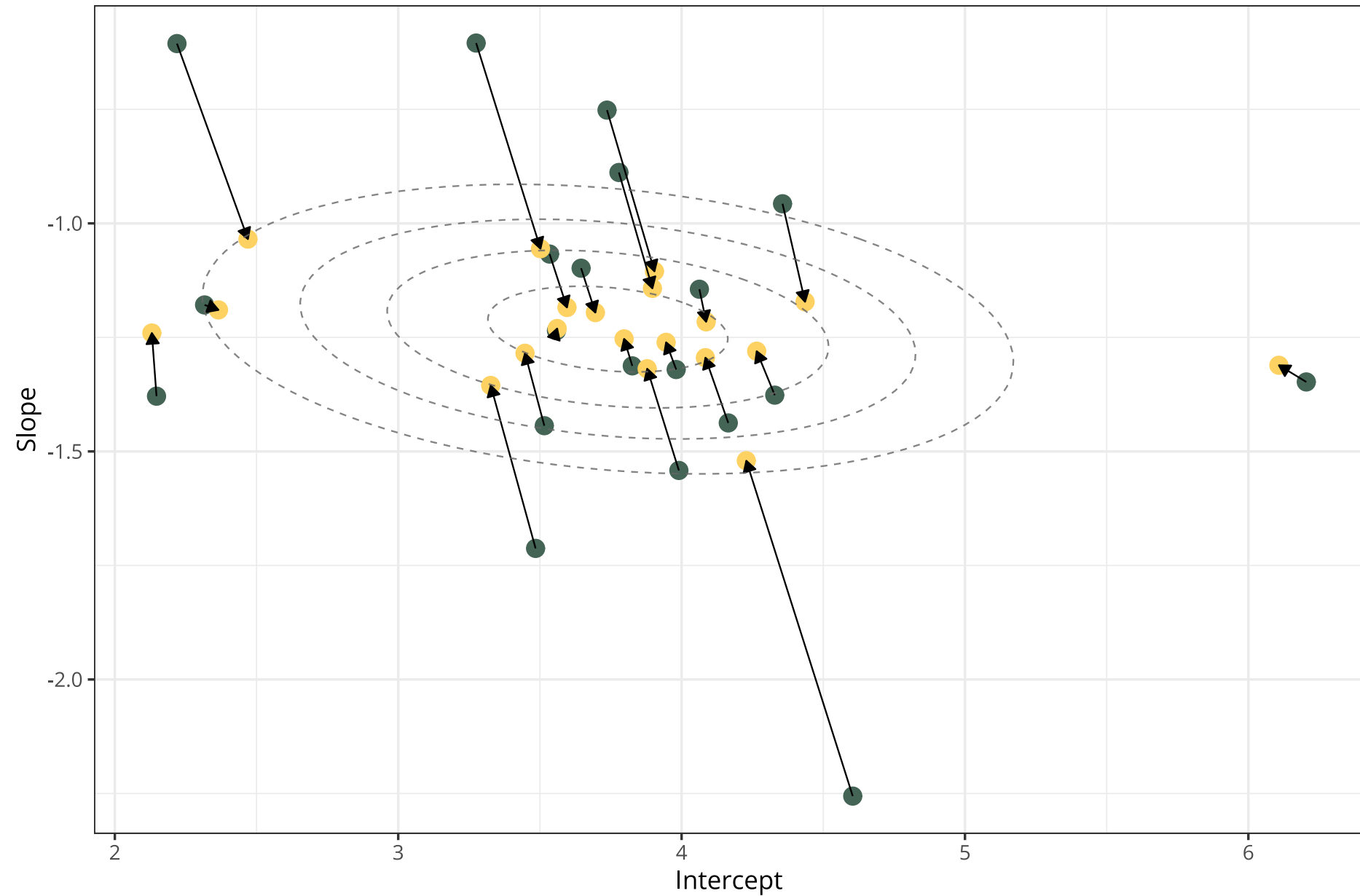
```

1 post <- as_draws_df(x = mod5) # extracts posterior samples
2 R <- rethinking::rlkjcorr(n = 16000, K = 2, eta = 2) # samples from prior
3
4 data.frame(prior = R[, 1, 2], posterior = post$cor_cafe__Intercept__afternoon) %>%
5   gather(type, value, prior:posterior) %>%
6   ggplot(aes(x = value, color = type, fill = type)) +
7   geom_histogram(position = "identity", alpha = 0.2) +
8   labs(x = expression(rho), y = "Nombre d'échantillons")

```



Shrinkage en deux dimensions



Comparaison de modèles

On compare le premier modèle (complete pooling model), le troisième modèle (partial pooling model), et le dernier modèle.

```
1 # comparaison des WAIC de chaque modèle
2 mod5 <- add_criterion(mod5, "waic")
3 w <- loo_compare(mod1, mod2, mod3, mod5, criterion = "waic")
4 print(w, simplify = FALSE)
```

	elpd_diff	se_diff	elpd_waic	se_elpd_waic	p_waic	se_p_waic	waic	se_waic
mod5	0.0	0.0	-155.0	10.1	26.6	2.6	310.0	20.1
mod3	-98.8	8.3	-253.8	8.3	18.1	1.5	507.6	16.6
mod2	-99.4	8.3	-254.4	8.4	19.4	1.6	508.7	16.7
mod1	-156.1	13.7	-311.1	10.6	2.0	0.3	622.1	21.1

```
1 model_weights(mod1, mod2, mod3, mod5, weights = "waic")
```

mod1	mod2	mod3	mod5
1.648031e-68	6.938094e-44	1.221044e-43	1.000000e+00



Comparaison de modèles

L'estimation du temps d'attente moyen est plus incertaine lorsqu'on prend en compte de nouvelles sources d'erreur. Cependant, l'erreur du modèle (i.e., ce qui n'est pas expliqué), la variation résiduelle σ , diminue...

```
1 posterior_summary(mod1, pars = c("^b", "sigma") )
```

	Estimate	Est.Error	Q2.5	Q97.5
b_Intercept	3.118824	0.07953888	2.963235	3.274122
sigma	1.141772	0.05642862	1.038975	1.256931

```
1 posterior_summary(mod3, pars = c("^b", "sigma") )
```

	Estimate	Est.Error	Q2.5	Q97.5
b_Intercept	3.1224462	0.20548659	2.7134453	3.5301519
sigma	0.8219151	0.04328802	0.7430558	0.9108058

```
1 posterior_summary(mod5, pars = c("^b", "sigma") )
```

	Estimate	Est.Error	Q2.5	Q97.5
b_Intercept	3.7397979	0.21646428	3.305611	4.1743387
b_afternoon	-1.2317902	0.08691248	-1.405042	-1.0591450
sigma	0.4892566	0.02748533	0.438913	0.5466079



Conclusions

Les modèles multi-niveaux (ou “modèles mixtes”) sont des extensions naturelles des modèles de régression classiques, où les paramètres de ces derniers se voient eux-même attribués des “modèles”, gouvernés par des hyper-paramètres.

Cette extension permet de faire des prédictions plus précises en prenant en compte la variabilité liée aux groupes ou structures (clusters) présent(e)s dans les données. Autrement dit, en modélisant les populations d'où sont tirés les effets aléatoires (e.g., la population de participants ou de stimuli).

Un modèle de régression classique est équivalent à un modèle multi-niveaux où la variabilité des effets aléatoires serait fixée à 0.

La cadre bayésien permet une interprétation naturelle des distributions desquelles proviennent les effets aléatoires (varying effects). En effet, ces distributions peuvent être interprétées comme des distributions a priori, dont les paramètres sont estimés à partir des données.



Travaux pratiques - sleepstudy

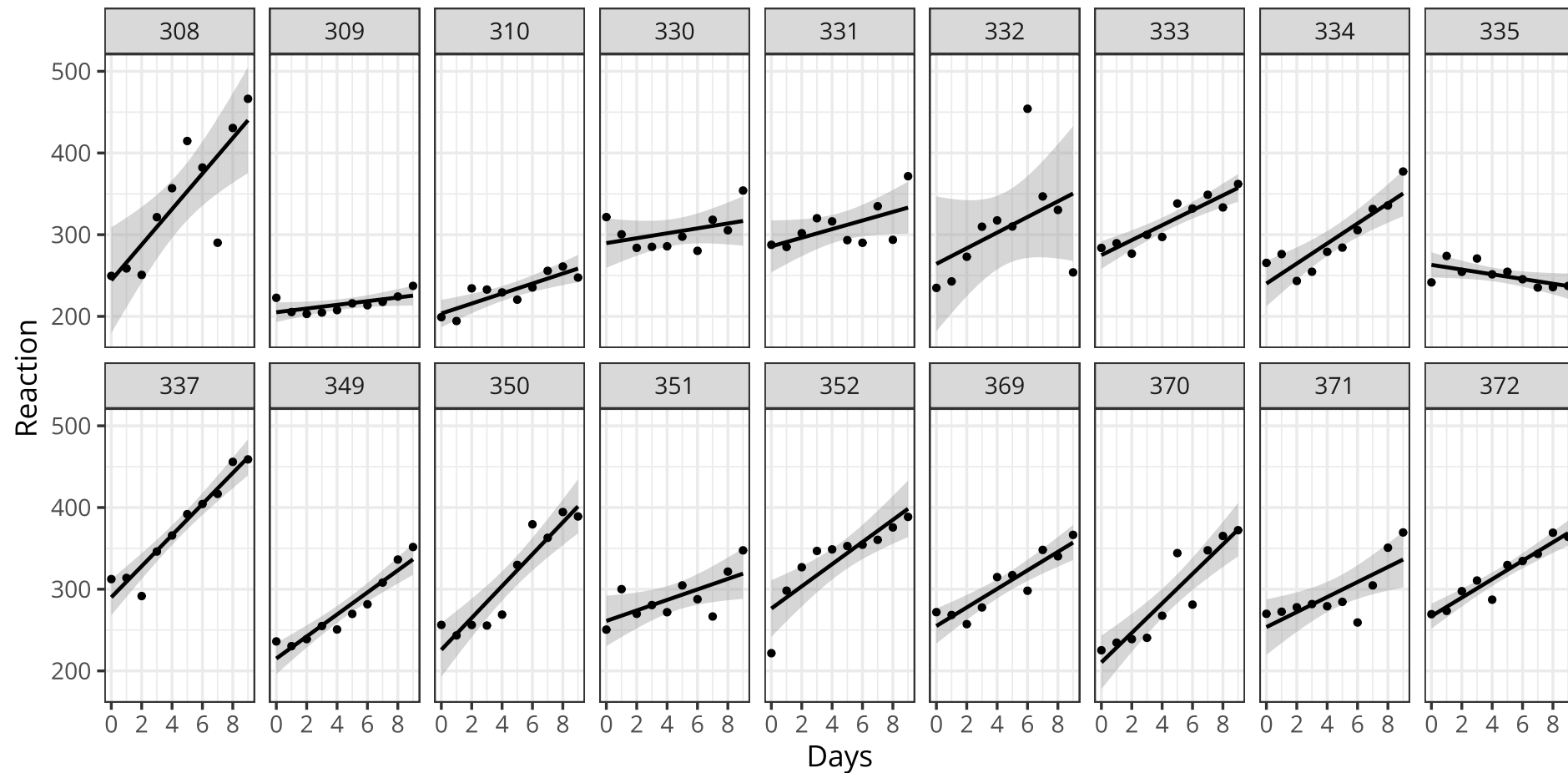
```
1 library(lme4)
2 data(sleepstudy)
3 head(sleepstudy, 20)
```

	Reaction	Days	Subject
1	249.5600	0	308
2	258.7047	1	308
3	250.8006	2	308
4	321.4398	3	308
5	356.8519	4	308
6	414.6901	5	308
7	382.2038	6	308
8	290.1486	7	308
9	430.5853	8	308
10	466.3535	9	308
11	222.7339	0	309
12	205.2658	1	309
13	202.9778	2	309
14	204.7070	3	309
15	207.7161	4	309
16	215.9618	5	309
17	213.6303	6	309
18	217.7272	7	309
19	224.2957	8	309
20	237.3142	9	309



Travaux pratiques - sleepstudy

```
1 sleepstudy %>%
2   ggplot(aes(x = Days, y = Reaction) ) +
3   geom_smooth(method = "lm", colour = "black") +
4   geom_point() +
5   facet_wrap(~Subject, nrow = 2) +
6   scale_x_continuous(breaks = c(0, 2, 4, 6, 8) )
```



Travaux pratiques - sleepstudy

À vous de construire les modèles mathématiques et les modèles **brms** correspondant aux modèles suivants :

- Modèle avec seulement l'effet fixe de **Days**.
- Modèle avec l'effet fixe de **Days** + un effet aléatoire de **Subject** (varying intercept).
- Modèle avec l'effet fixe de **Days** + un effet aléatoire de **Subject**. (varying intercept) + un effet aléatoire de **Days** (varying slope).

Comparez ensuite ces modèles en utilisant les outils discutés aux cours précédents (e.g., WAIC) et concluez.



Proposition de solution

```
1 fmod0 <- lm(Reaction ~ Days, sleepstudy)
2 fmod1 <- lmer(Reaction ~ Days + (1 | Subject), sleepstudy)
3 fmod2 <- lmer(Reaction ~ Days + (1 + Days | Subject), sleepstudy)
4
5 anova(fmod1, fmod2)
```

Data: sleepstudy

Models:

fmod1: Reaction ~ Days + (1 | Subject)

fmod2: Reaction ~ Days + (1 + Days | Subject)

	npar	AIC	BIC	logLik	deviance	Chisq	Df	Pr(>Chisq)
fmod1	4	1802.1	1814.8	-897.04	1794.1			
fmod2	6	1763.9	1783.1	-875.97	1751.9	42.139	2	7.072e-10 ***

Signif. codes: 0 '***' 0.001 '**' 0.01 '*' 0.05 '.' 0.1 ' ' 1



Proposition de solution

```

1 mod6 <- brm(
2   Reaction ~ 1 + Days,
3   prior = c(
4     prior(normal(200, 100), class = Intercept),
5     prior(normal(0, 10), class = b),
6     prior(cauchy(0, 10), class = sigma)
7   ),
8   data = sleepstudy,
9   warmup = 1000, iter = 5000,
10  cores = parallel::detectCores()
11  )

```

```
1 posterior_summary(mod6)
```

	Estimate	Est.Error	Q2.5	Q97.5
b_Intercept	251.93039	6.6371239	238.884234	264.84498
b_Days	10.32424	1.2315772	7.902365	12.73023
sigma	47.74154	2.5503911	43.037914	52.96216
lprior	-15.69187	0.1670769	-16.039829	-15.38754
lp__	-963.48096	1.2294009	-966.624701	-962.07437



Proposition de solution

```

1 mod7 <- brm(
2   Reaction ~ 1 + Days + (1 | Subject),
3   prior = c(
4     prior(normal(200, 100), class = Intercept),
5     prior(normal(0, 10), class = b),
6     prior(cauchy(0, 10), class = sigma),
7     prior(cauchy(0, 10), class = sd)
8   ),
9   data = sleepstudy,
10  warmup = 1000, iter = 5000,
11  cores = parallel::detectCores()
12 )

```

```

1 posterior_summary(mod7, pars = c("^b", "sigma") )

```

	Estimate	Est.Error	Q2.5	Q97.5
b_Intercept	250.88978	10.0223168	230.99318	270.56323
b_Days	10.40515	0.8027281	8.84247	11.99794
sigma	31.08754	1.7404375	27.86194	34.74327



Proposition de solution

```

1 mod8 <- brm(
2   Reaction ~ 1 + Days + (1 + Days | Subject),
3   prior = c(
4     prior(normal(200, 100), class = Intercept),
5     prior(normal(0, 10), class = b),
6     prior(cauchy(0, 10), class = sigma),
7     prior(cauchy(0, 10), class = sd)
8   ),
9   data = sleepstudy,
10  warmup = 1000, iter = 5000,
11  cores = parallel::detectCores()
12 )

```

```

1 posterior_summary(mod8, pars = c("^b", "sigma") )

```

	Estimate	Est.Error	Q2.5	Q97.5
b_Intercept	251.14710	6.934251	237.643148	264.81612
b_Days	10.07880	1.635332	6.724357	13.21770
sigma	25.85959	1.555886	23.008627	29.06938



Proposition de solution

```

1 # calcul du WAIC et ajout du WAIC à chaque modèle
2 mod6 <- add_criterion(mod6, "waic")
3 mod7 <- add_criterion(mod7, "waic")
4 mod8 <- add_criterion(mod8, "waic")
5
6 # comparaison des WAIC de chaque modèle
7 w <- loo_compare(mod6, mod7, mod8, criterion = "waic")
8 print(w, simplify = FALSE)

```

	elpd_diff	se_diff	elpd_waic	se_elpd_waic	p_waic	se_p_waic	waic	se_waic
mod8	0.0	0.0	-860.2	22.3	32.7	8.3	1720.5	44.6
mod7	-24.4	11.5	-884.6	14.4	19.1	3.3	1769.2	28.8
mod6	-93.1	20.9	-953.3	10.6	3.2	0.5	1906.6	21.2

```

1 # calcul du poids relatif de chaque modèle
2 model_weights(mod6, mod7, mod8, weights = "waic")

```

mod6	mod7	mod8
3.798293e-41	2.647074e-11	1.000000e+00



Références

Efron, B., & Morris, C. (1977). Stein's paradox in statistics. *Scientific American*, 236(5), 119–127.

<https://doi.org/10.1038/scientificamerican0577-119>

Gelman, A. (2005). Analysis of variance? Why it is more important than ever. *The Annals of Statistics*, 33(1), 1–53. <https://doi.org/10.1214/009053604000001048>

Gelman, A., & Hill, J. (2006). *Data analysis using regression and multilevel/hierarchical models*.
<https://doi.org/10.1017/cbo9780511790942>

Lewandowski, D., Kurowicka, D., & Joe, H. (2009). Generating random correlation matrices based on vines and extended onion method. *Journal of Multivariate Analysis*, 100(9), 1989–2001.

<https://doi.org/10.1016/j.jmva.2009.04.008>

