

Introduction à la modélisation statistique bayésienne

Un cours en R, Stan, et 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



Null Hypothesis Significance Testing (NHST)

On s'intéresse aux différences de taille entre hommes et femmes. On va mesurer 100 femmes et 100 hommes.

```
1 set.seed(19) # pour reproduire les résultats
2 men <- rnorm(100, 175, 10) # 100 tailles d'hommes
3 women <- rnorm(100, 170, 10) # 100 tailles de femmes
```

```
1 t.test(men, women) # test de student pour différence de moyenne
```

```
Welch Two Sample t-test

data:  men and women
t = 2.4969, df = 197.98, p-value = 0.01335
alternative hypothesis: true difference in means is not equal to 0
95 percent confidence interval:
 0.7658541 6.5209105
sample estimates:
mean of x mean of y
175.1258 171.4825
```





Null Hypothesis Significance Testing (NHST)

On va simuler des t-valeurs issues de données générées sous l'hypothèse d'une absence de différence entre hommes et femmes.

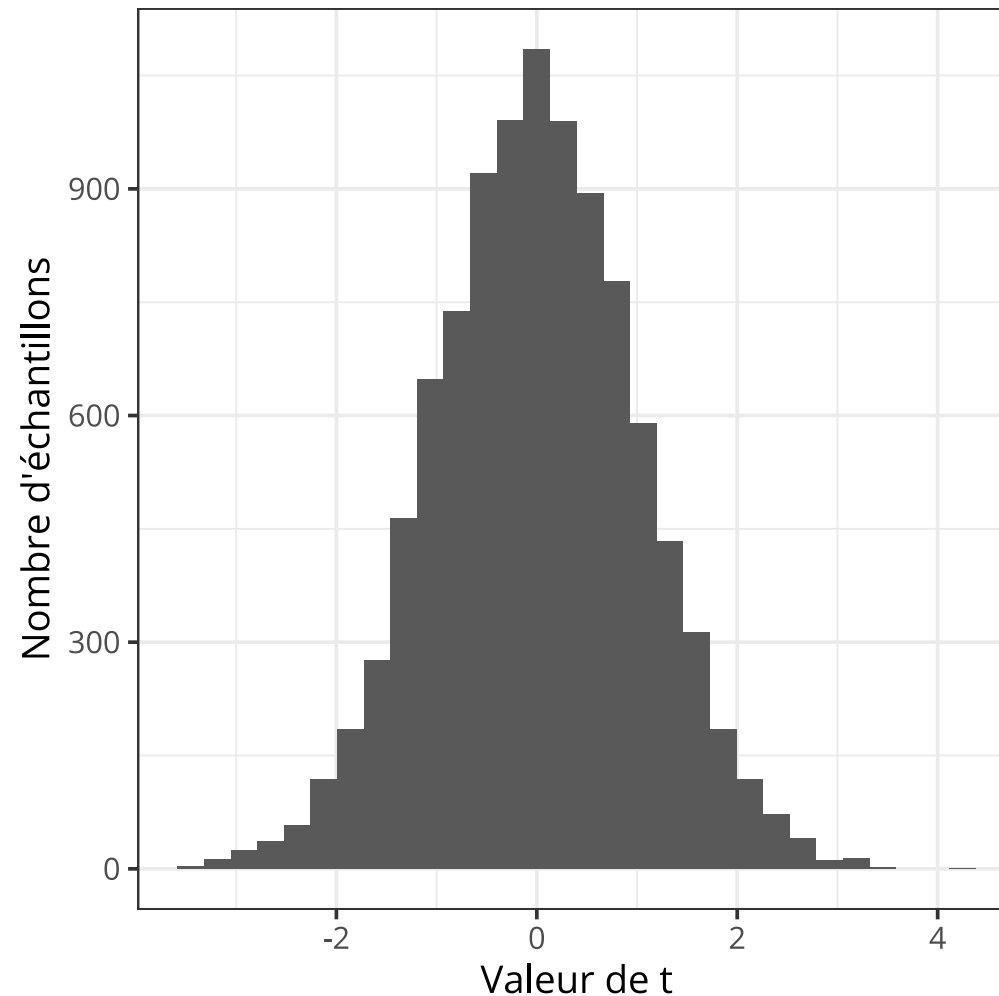
```
1 nsims <- 1e4 # nombre de simulations
2 t <- rep(NA, nsims) # initialisation d'un vecteur vide
3
4 for (i in 1:nsims) {
5
6   men2 <- rnorm(100, 170, 10) # 100 tailles d'hommes
7   women2 <- rnorm(100, 170, 10) # 100 tailles de femmes (de même distribution)
8   t[i] <- t.test(men2, women2)$statistic # on conserve la t-valeur
9
10 }
```

```
1 # une autre manière de réaliser la même opération, sans boucle for
2 t <- replicate(nsims, t.test(rnorm(100, 170, 10), rnorm(100, 170, 10))$statistic)
```



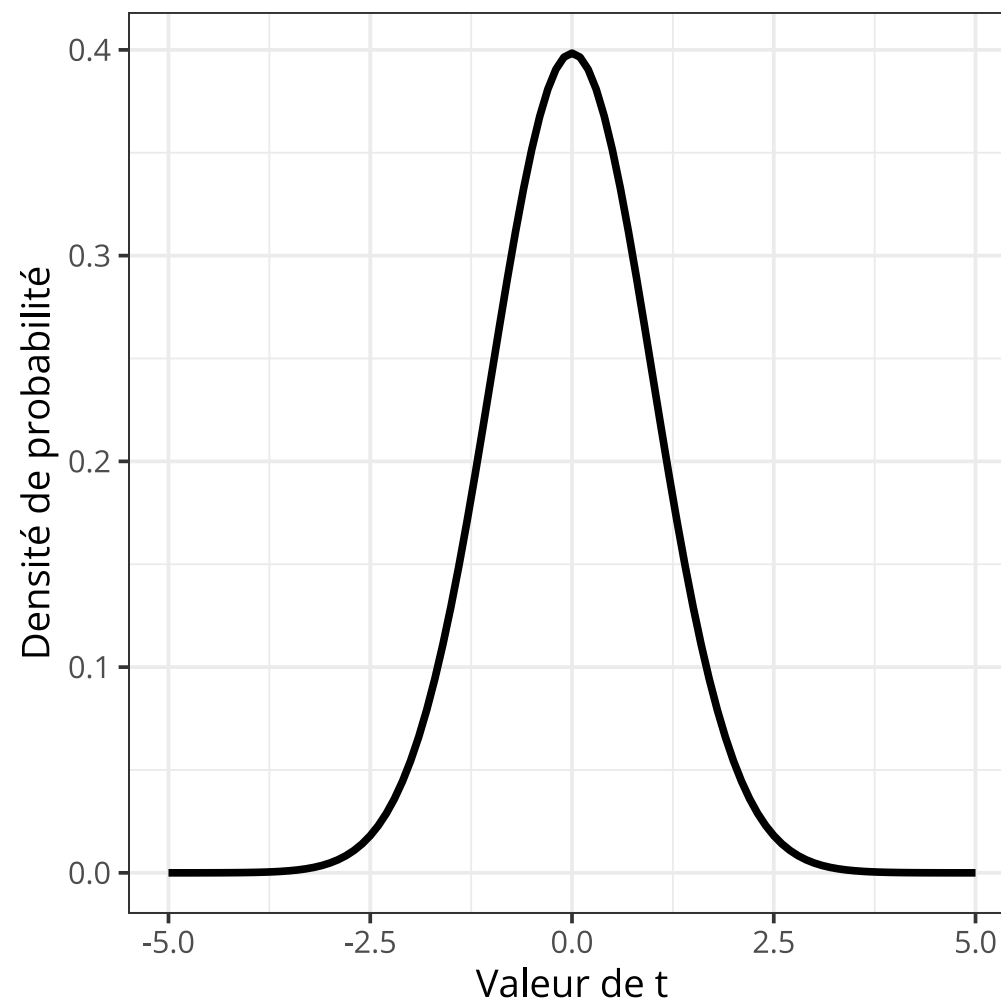
Null Hypothesis Significance Testing (NHST)

```
1 data.frame(t = t) %>%  
2   ggplot(aes(x = t) ) +  
3   geom_histogram() +  
4   labs(x = "Valeur de t", y = "Nombre d'échantillons")
```



Null Hypothesis Significance Testing (NHST)

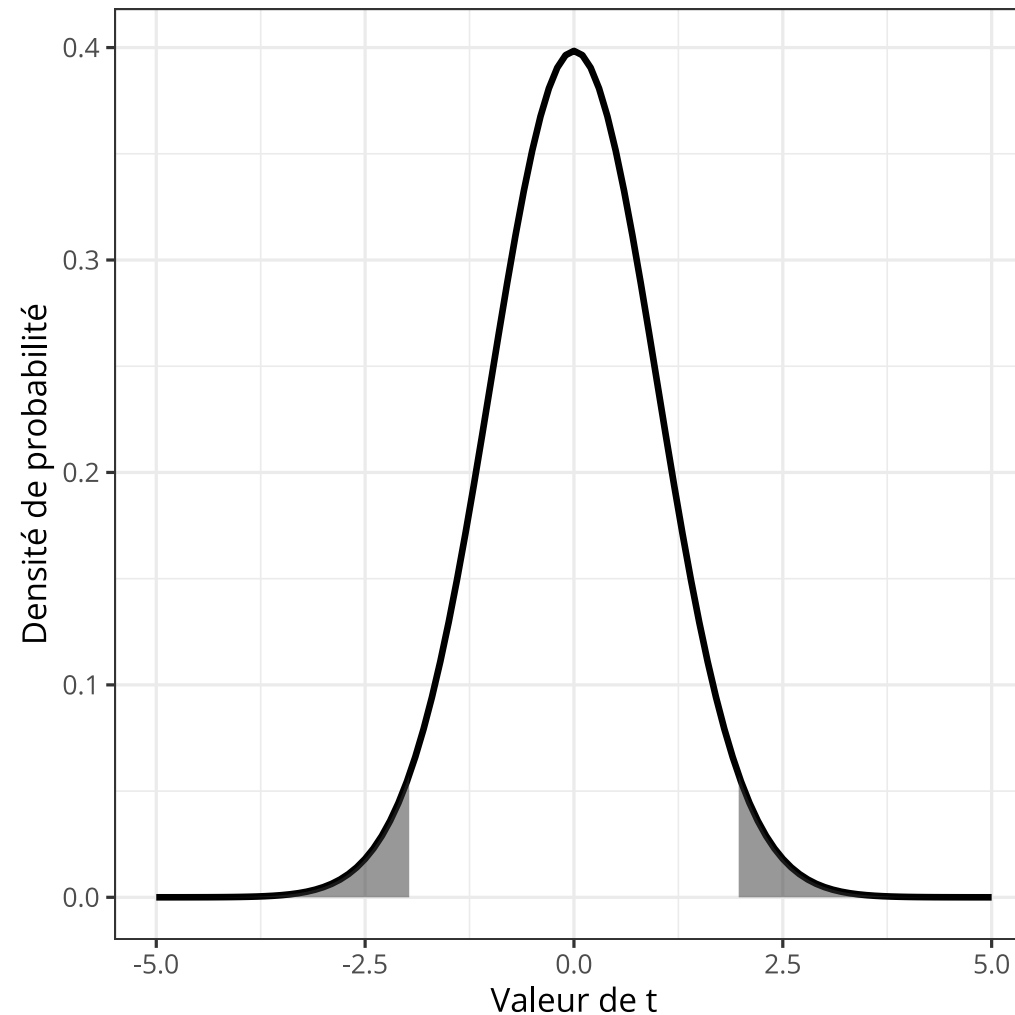
```
1 data.frame(x = c(-5, 5) ) %>%  
2   ggplot(aes(x = x) ) +  
3   stat_function(fun = dt, args = list(df = t.test(men, women)$parameter), size = 1.5) +  
4   labs(x = "Valeur de t", y = "Densité de probabilité")
```



Null Hypothesis Significance Testing (NHST)

```
1 alpha <- 0.05 # seuil de significativité  
2 abs(qt(alpha / 2, df = t.test(men, women)$parameter) ) # valeur de t critique
```

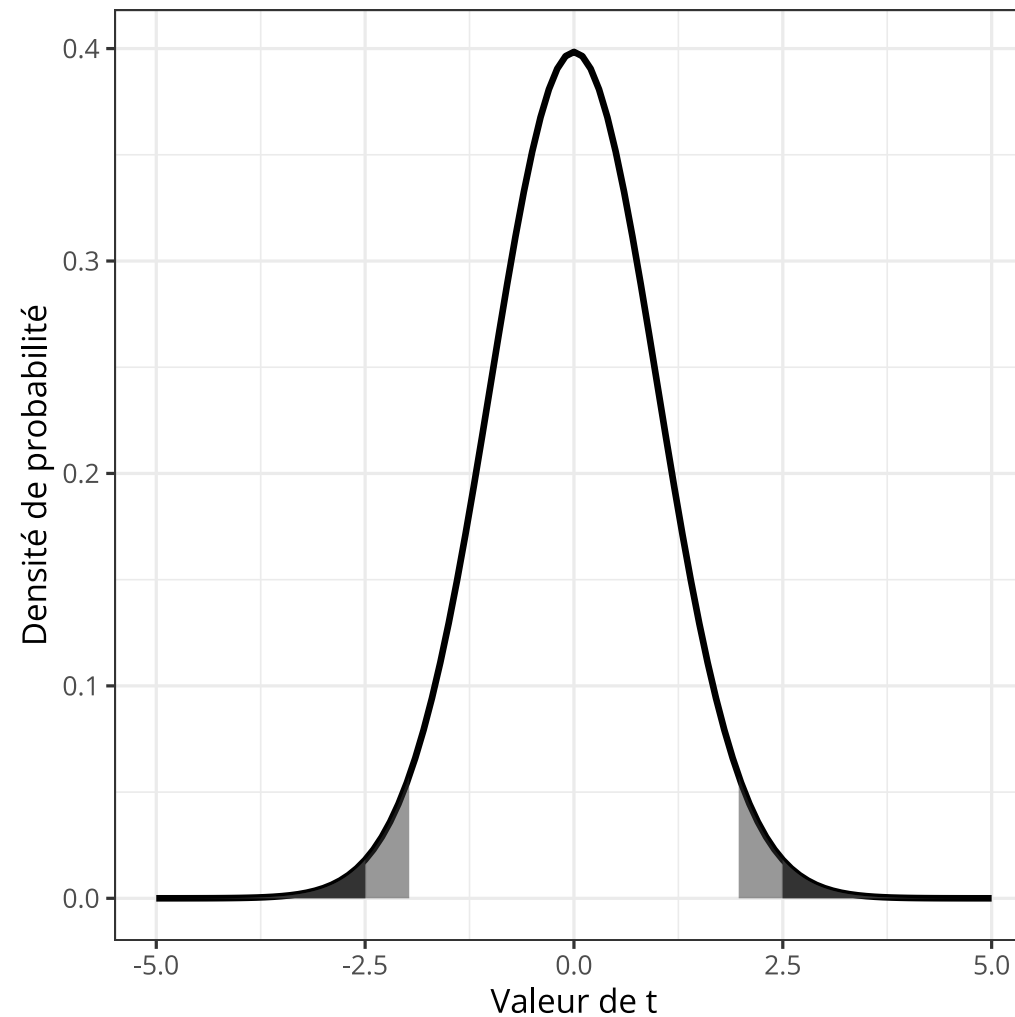
```
[1] 1.972019
```



Null Hypothesis Significance Testing (NHST)

```
1 tobs <- t.test(men, women)$statistic # valeur de t observée  
2 tobs %>% as.numeric
```

```
[1] 2.496871
```



P-valeur

Une p-valeur est une aire sous la courbe (une intégrale) sous la distribution de statistiques de test sous l'hypothèse nulle (i.e., étant admis que l'hypothèse nulle est vraie). La p-valeur indique la probabilité d'observer la valeur de la statistique de test (e.g., une valeur de t) observée, ou une valeur plus extrême, sous l'hypothèse nulle.

$$p[t(\mathbf{x}^{\text{rep}}; \mathcal{H}_0) \geq t(x)]$$

```
1 t.test(men, women)$p.value
```

```
[1] 0.01334509
```

```
1 tvalue <- abs(t.test(men, women)$statistic)
2 df <- t.test(men, women)$parameter
3 2 * integrate(dt, tvalue, Inf, df = df)$value
```

```
[1] 0.01334509
```

```
1 2 * (1 - pt(abs(t.test(men, women)$statistic), t.test(men, women)$parameter) )
```

```
      t
0.01334509
```



Intervalles de confiance

Les intervalles de confiance peuvent être interprétés comme des régions de significativité (ou des régions de **compatibilité** avec l'hypothèse nulle). Par conséquent, un intervalle de confiance contient la même information qu'une p-valeur et s'interprète de manière similaire.

On ne peut pas dire qu'un intervalle de confiance a une probabilité de 95% de contenir la vraie valeur (i.e., la valeur dans la population) de θ (cf. the [inverse fallacy](#)), contrairement à l'intervalle de crédibilité bayésien.

Un intervalle de confiance à 95% représente un degré de “recouvrement” (coverage). Le “95%” fait référence à une propriété fréquentiste (i.e., sur le long-terme) de la procédure, mais ne fait pas référence au paramètre θ . Autrement dit, sur le long-terme, 95% des intervalles de confiance à 95% que l'on pourrait calculer (dans une réplique exacte de notre expérience) contiendraient la valeur du paramètre dans la population (i.e., la “vraie” valeur de θ). Cependant, nous ne pouvons pas dire qu'un intervalle de confiance en particulier a une probabilité de 95% de contenir la “vraie” valeur de θ ... soit ce dernier contient la vraie valeur de θ , soit il ne la contient pas.





Facteur de Bayes

On compare deux modèles :

- $\mathcal{H}_0 : \mu_1 = \mu_2 \rightarrow \delta = 0$
- $\mathcal{H}_1 : \mu_1 \neq \mu_2 \rightarrow \delta \neq 0$

$$\underbrace{\frac{p(\mathcal{H}_0 | D)}{p(\mathcal{H}_1 | D)}}_{\text{posterior odds}} = \underbrace{\frac{p(D | \mathcal{H}_0)}{p(D | \mathcal{H}_1)}}_{\text{Bayes factor}} \times \underbrace{\frac{p(\mathcal{H}_0)}{p(\mathcal{H}_1)}}_{\text{prior odds}}$$

$$\text{evidence} = p(D | \mathcal{H}) = \int p(\theta | \mathcal{H})p(D | \theta, \mathcal{H})d\theta$$

L'évidence en faveur d'un modèle correspond à la **vraisemblance marginale** d'un modèle (le dénominateur du théorème de Bayes), c'est à dire à la vraisemblance moyennée sur le prior... Ce qui fait du facteur de Bayes un ratio de vraisemblances, pondéré par (ou moyenné sur) le prior.

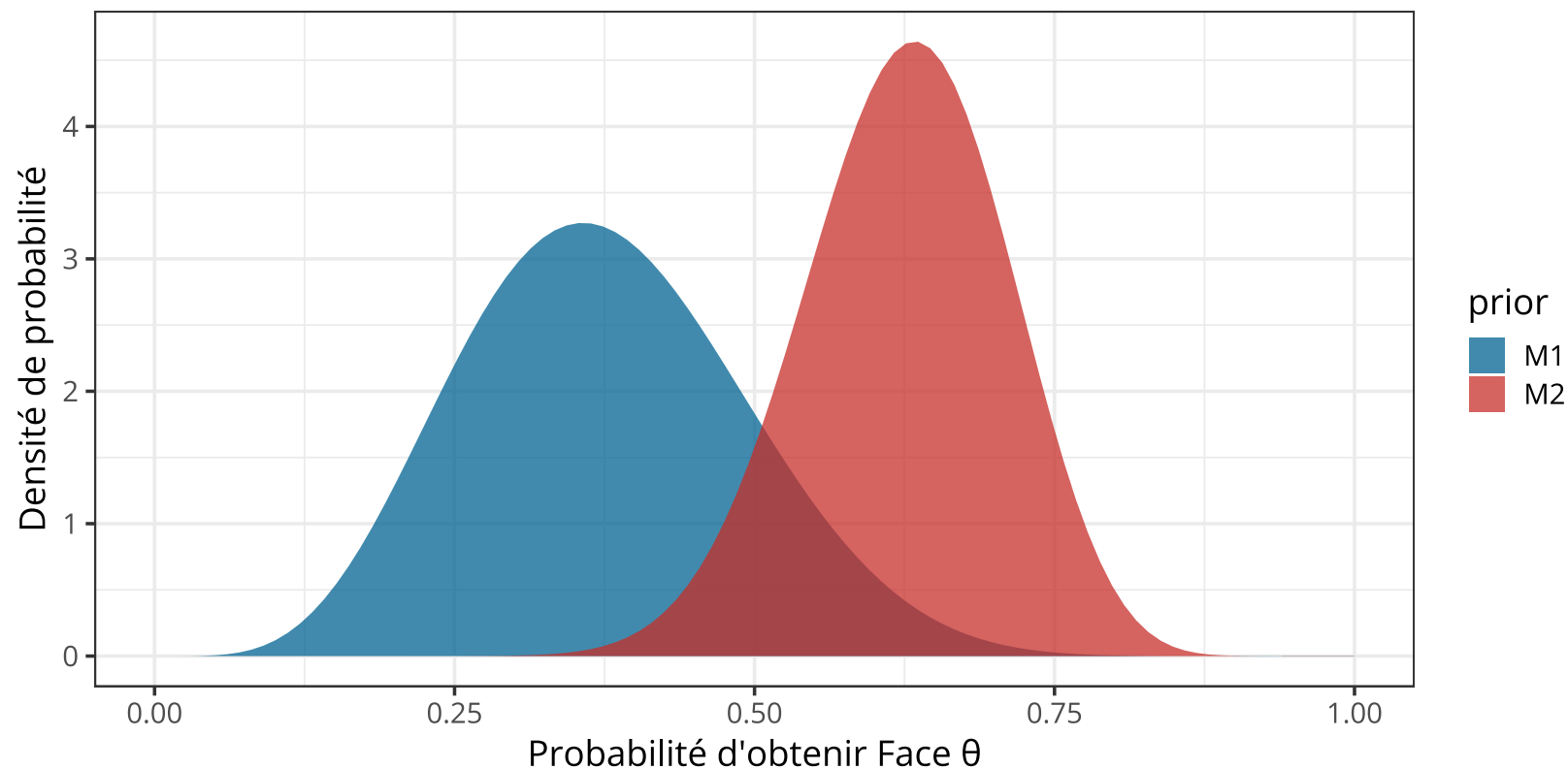


Facteur de Bayes, exemple d'application

On lance une pièce 100 fois et on essaye d'estimer la probabilité θ (le biais de la pièce) d'obtenir Face. On compare deux modèles qui diffèrent par leur prior sur θ .

$$\mathcal{M}_1 : y_i \sim \text{Binomial}(n, \theta) \\ \theta \sim \text{Beta}(6, 10)$$

$$\mathcal{M}_2 : y_i \sim \text{Binomial}(n, \theta) \\ \theta \sim \text{Beta}(20, 12)$$



Facteur de Bayes, exemple d'application

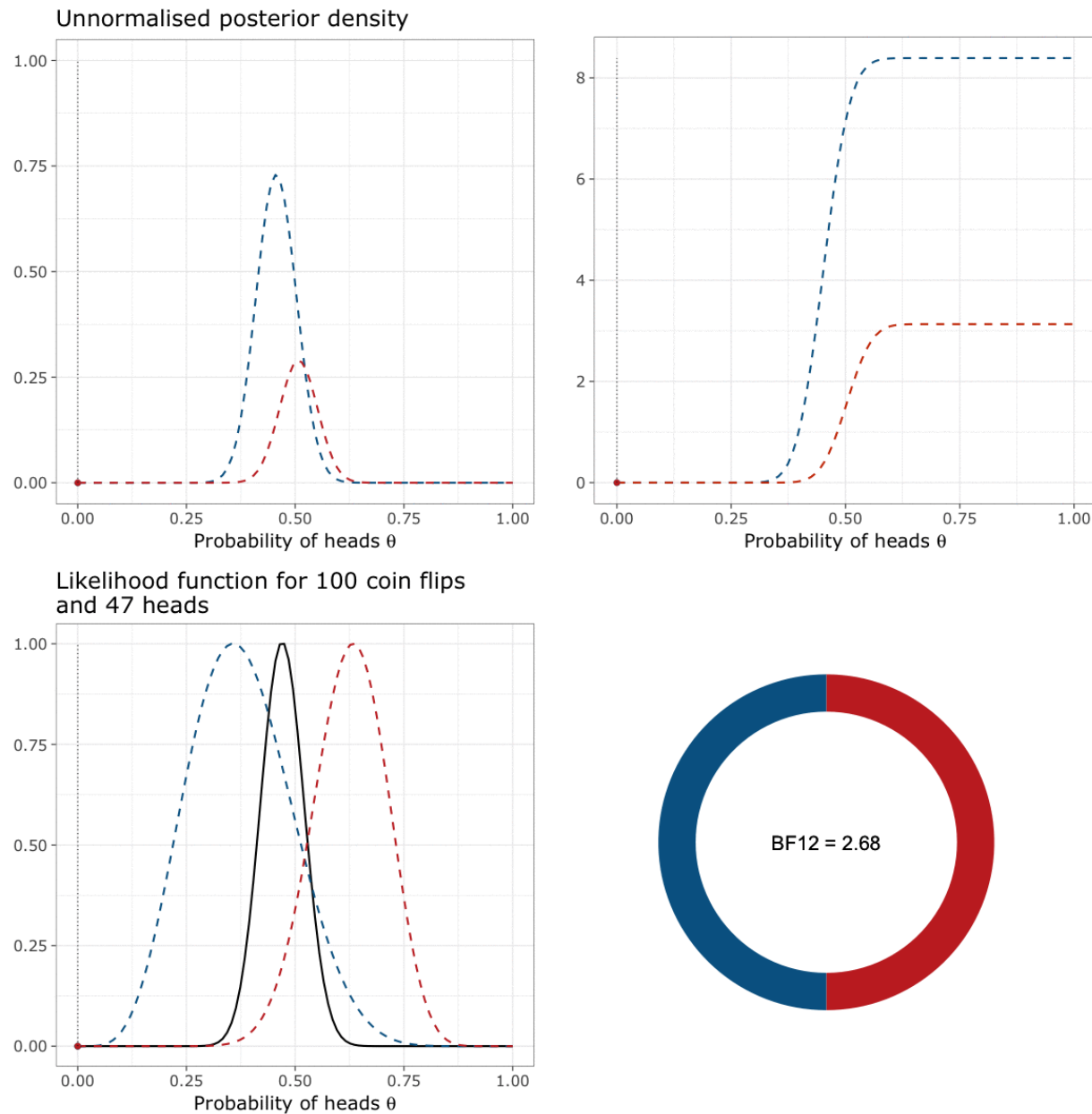
$$\begin{aligned}\mathcal{M}_1 : y_i &\sim \text{Binomial}(n, \theta) \\ \theta &\sim \text{Beta}(6, 10)\end{aligned}$$

$$\begin{aligned}\mathcal{M}_2 : y_i &\sim \text{Binomial}(n, \theta) \\ \theta &\sim \text{Beta}(20, 12)\end{aligned}$$

$$\text{BF}_{12} = \frac{p(D|\mathcal{M}_1)}{p(D|\mathcal{M}_2)} = \frac{\int p(\theta|\mathcal{M}_1)p(D|\theta, \mathcal{M}_1)d\theta}{\int p(\theta|\mathcal{M}_2)p(D|\theta, \mathcal{M}_2)d\theta} = \frac{\int \text{Binomial}(n, \theta)\text{Beta}(6, 10)d\theta}{\int \text{Binomial}(n, \theta)\text{Beta}(20, 12)d\theta}$$



Facteur de Bayes, exemple d'application



Le facteur de Bayes est la nouvelle p-valeur

Attention à ne pas interpréter le BF comme un rapport des chances a posteriori (posterior odds)...

Le BF est un facteur qui nous indique de combien notre rapport des chances a priori (prior odds) doit changer, au vu des données. Il ne nous dit pas quelle est l'hypothèse la plus probable, sachant les données (sauf si les prior odds sont de 1/1). Par exemple :

- \mathcal{H}_0 : la précognition n'existe pas !
- \mathcal{H}_1 : la précognition existe.

On fait une expérience et on calcule un $\text{BF}_{10} = 27$. Quelle est la probabilité a posteriori de \mathcal{H}_1 ?

$$\underbrace{\frac{p(\mathcal{H}_1 | D)}{p(\mathcal{H}_0 | D)}}_{\text{posterior odds}} = \underbrace{\frac{27}{1}}_{\text{Bayes factor}} \times \underbrace{\frac{1}{1000}}_{\text{prior odds}} = \frac{27}{1000} = 0.027$$



Double sens

Dans le cadre bayésien, le terme de **prior** peut faire référence soit :

- À la probabilité a priori ou a posteriori d'un modèle (par rapport à un autre modèle), c'est à dire $\Pr(\mathcal{M}_i)$. Voir ce [blogpost](#).
- Aux priors qu'on définit sur les paramètres d'un modèle, par exemple $\beta \sim \text{Normal}(0, 1)$. Voir ce [blogpost](#).



Comparaison de modèles

Deux problèmes récurrents à éviter en modélisation : le sous-apprentissage (underfitting) et sur-apprentissage (overfitting). Comment s'en sortir ?

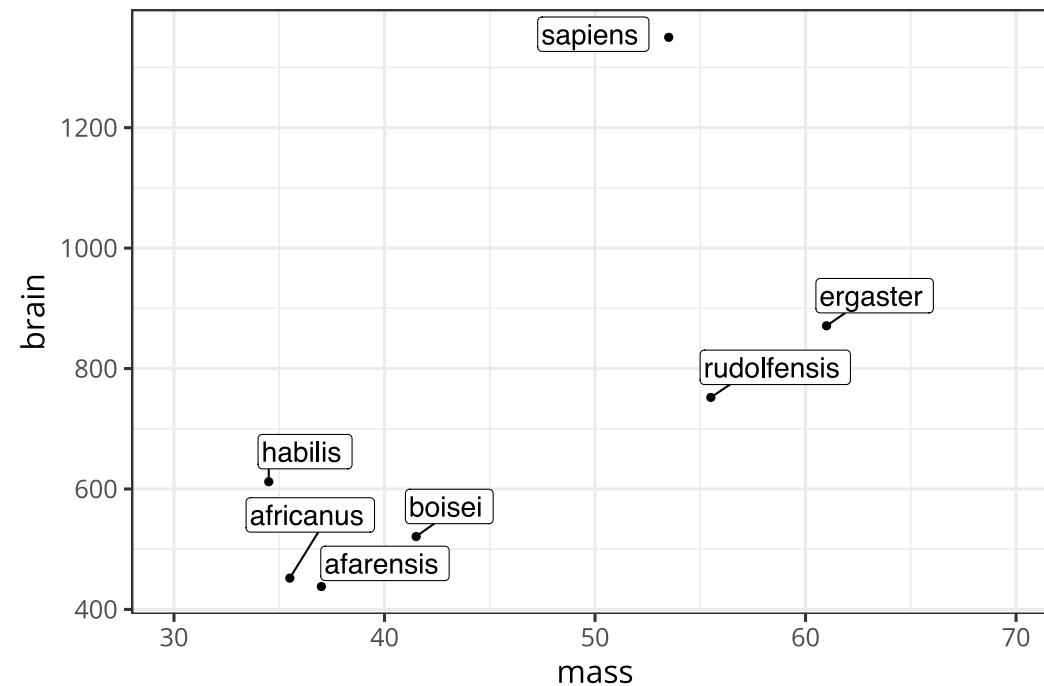
- Utiliser des priors pour **régulariser**, pour contraindre le posterior (i.e., accorder moins de poids à la vraisemblance).
- Utiliser des critères d'information (e.g., AIC, WAIC).

$$R^2 = \frac{\text{var}(\text{outcome}) - \text{var}(\text{residuals})}{\text{var}(\text{outcome})} = 1 - \frac{\text{var}(\text{residuals})}{\text{var}(\text{outcome})}$$



Overfitting

```
1 ppnames <- c("afarensis", "africanus", "habilis", "boisei",  
2             "rudolfensis", "ergaster", "sapiens")  
3 brainvolcc <- c(438, 452, 612, 521, 752, 871, 1350)  
4 masskg <- c(37.0, 35.5, 34.5, 41.5, 55.5, 61.0, 53.5)  
5  
6 d <- data.frame(species = ppnames, brain = brainvolcc, mass = masskg)  
7  
8 d %>%  
9   ggplot(aes(x = mass, y = brain, label = species)) +  
10  geom_point() +  
11  ggrepel::geom_label_repel(hjust = 0, nudge_y = 50, size = 5) +  
12  xlim(30, 70)
```



Overfitting

```
1 mod1.1 <- lm(brain ~ mass, data = d)
2 (var(d$brain) - var(residuals(mod1.1)) ) / var(d$brain)
```

```
[1] 0.490158
```

```
1 mod1.2 <- lm(brain ~ mass + I(mass^2), data = d)
2 (var(d$brain) - var(residuals(mod1.2)) ) / var(d$brain)
```

```
[1] 0.5359967
```

```
1 mod1.3 <- lm(brain ~ mass + I(mass^2) + I(mass^3), data = d)
2 (var(d$brain) - var(residuals(mod1.3)) ) / var(d$brain)
```

```
[1] 0.6797736
```



Overfitting

```
1 mod1.4 <- lm(brain ~ mass + I(mass^2) + I(mass^3) + I(mass^4), data = d)
2 (var(d$brain) - var(residuals(mod1.4)) ) / var(d$brain)
```

```
[1] 0.8144339
```

```
1 mod1.5 <- lm(brain ~ mass + I(mass^2) + I(mass^3) + I(mass^4) +
2             I(mass^5), data = d)
3 (var(d$brain) - var(residuals(mod1.5)) ) / var(d$brain)
```

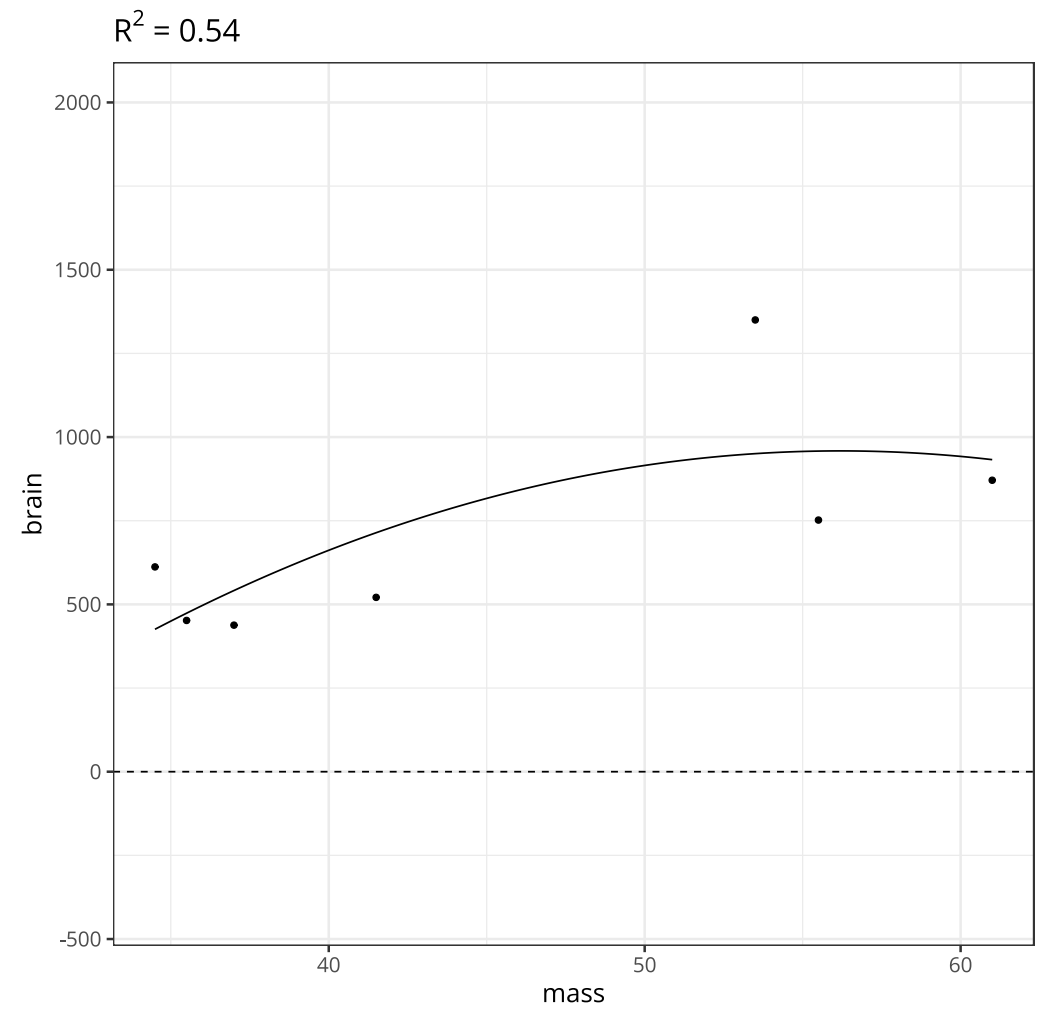
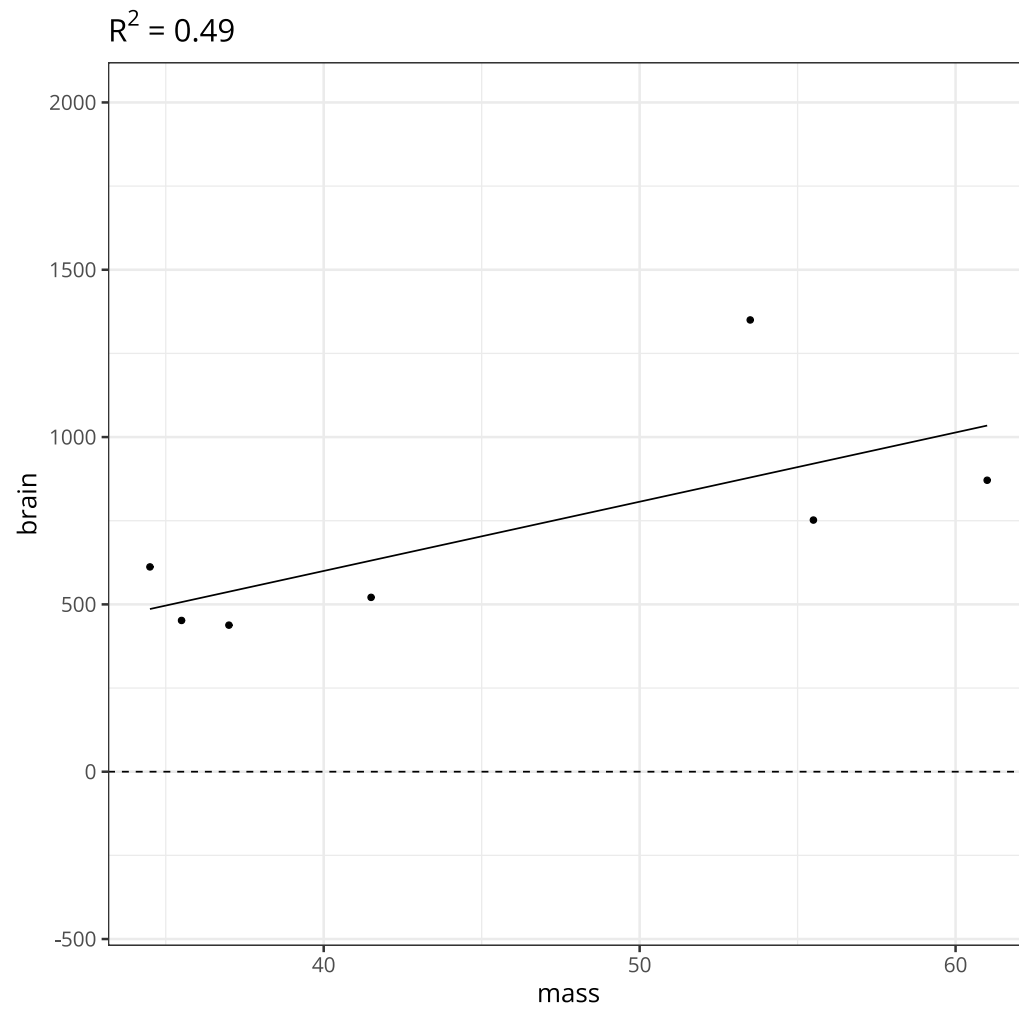
```
[1] 0.988854
```

```
1 mod1.6 <- lm( brain ~ mass + I(mass^2) + I(mass^3) + I(mass^4) +
2             I(mass^5) + I(mass^6), data = d)
3 (var(d$brain) - var(residuals(mod1.6)) ) / var(d$brain)
```

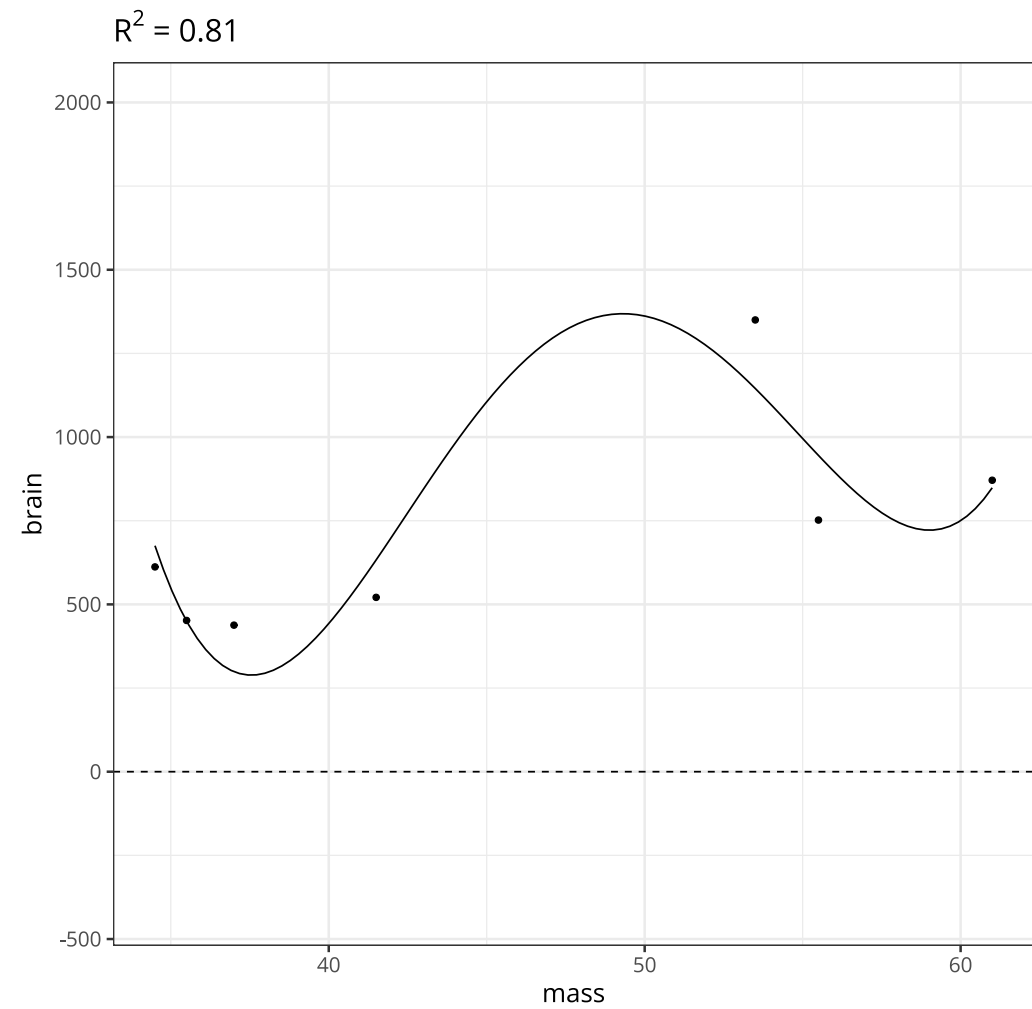
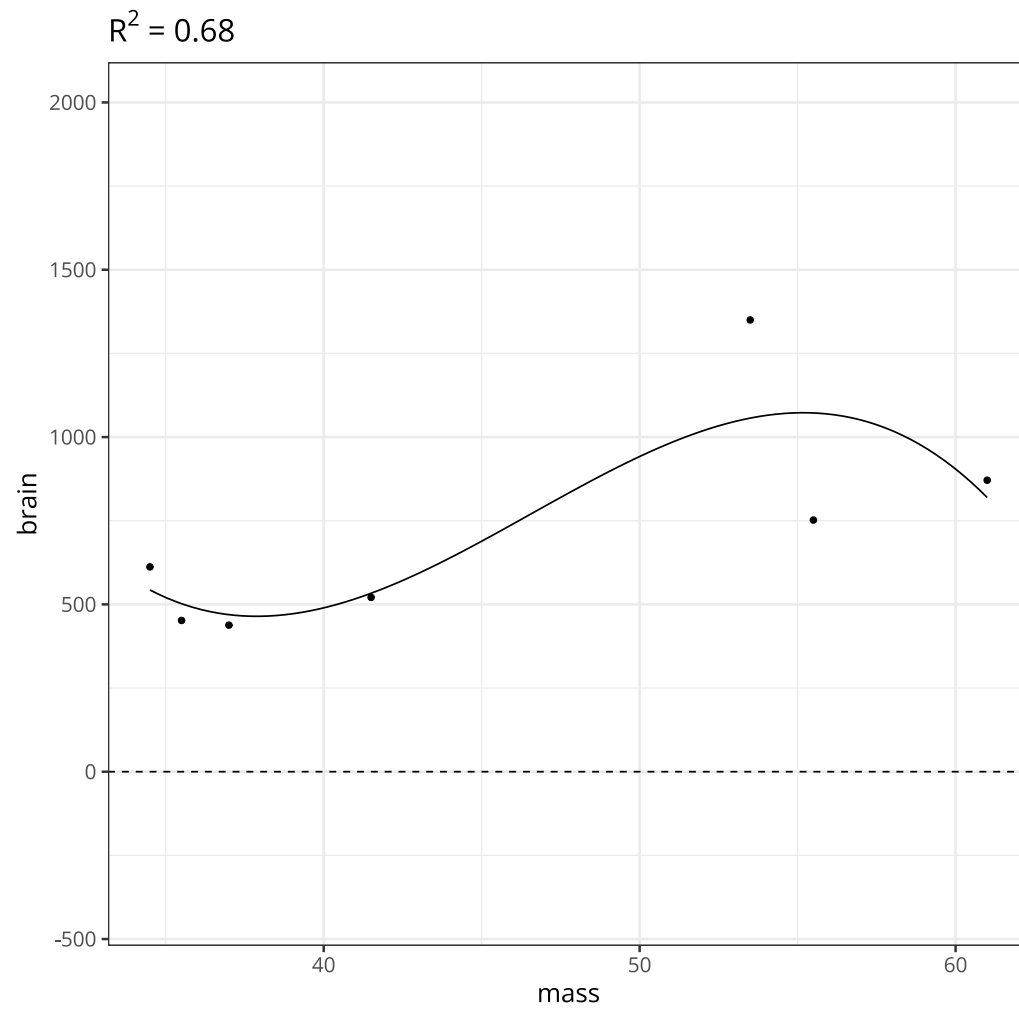
```
[1] 1
```



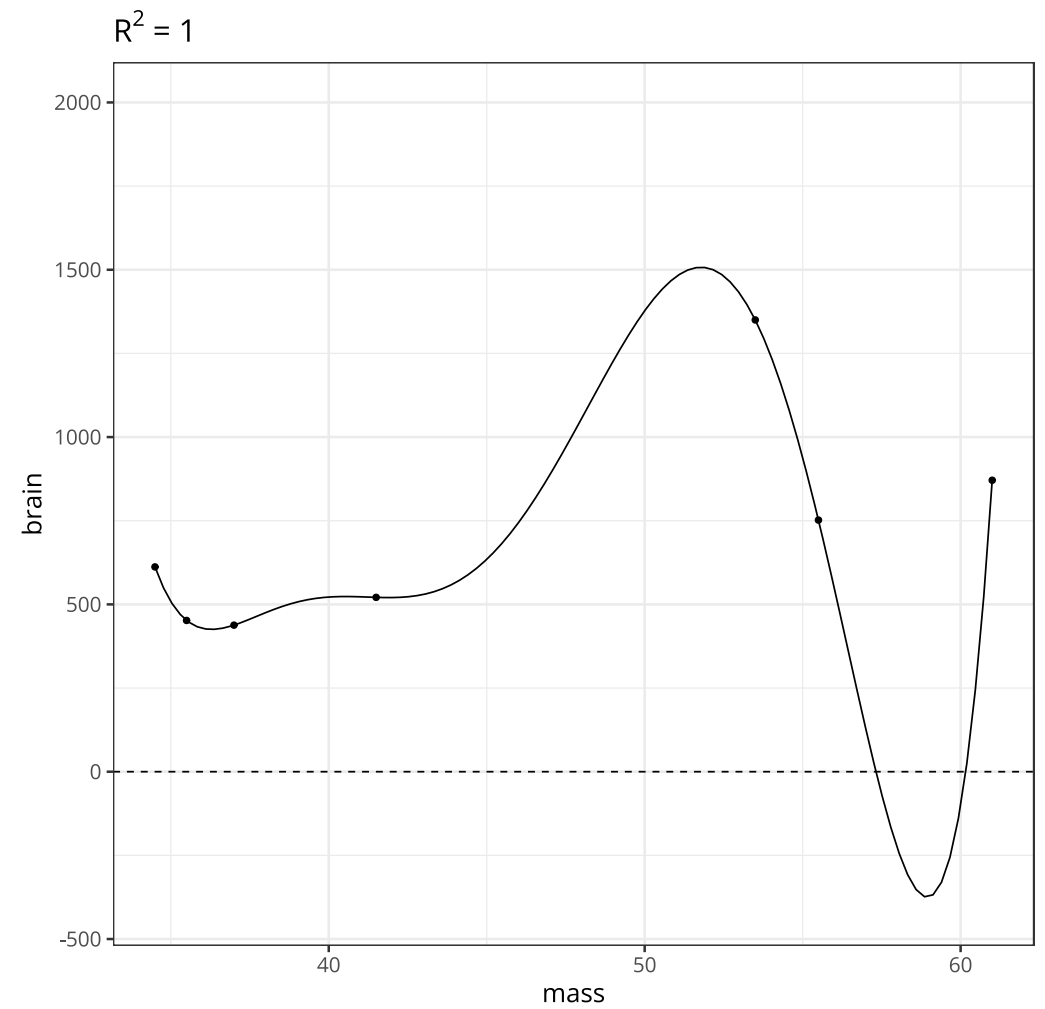
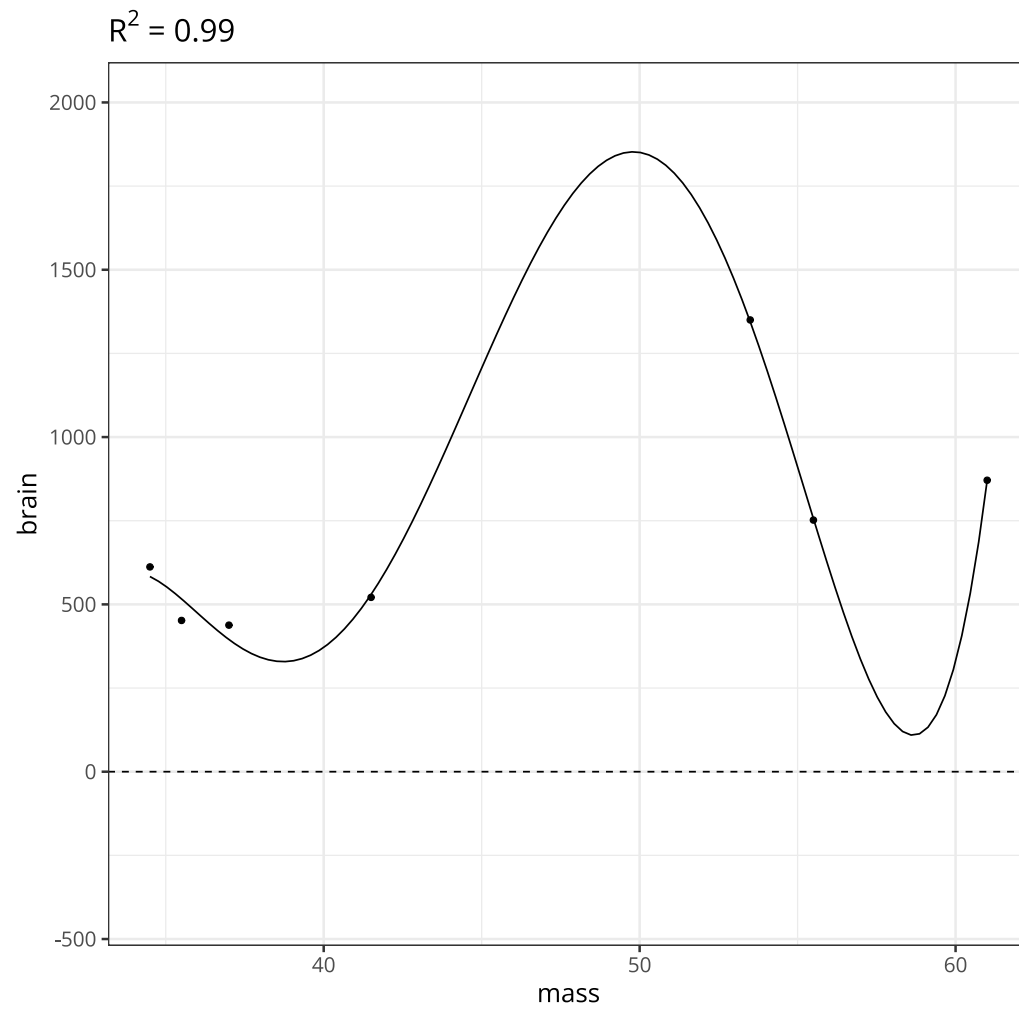
Overfitting



Overfitting



Overfitting

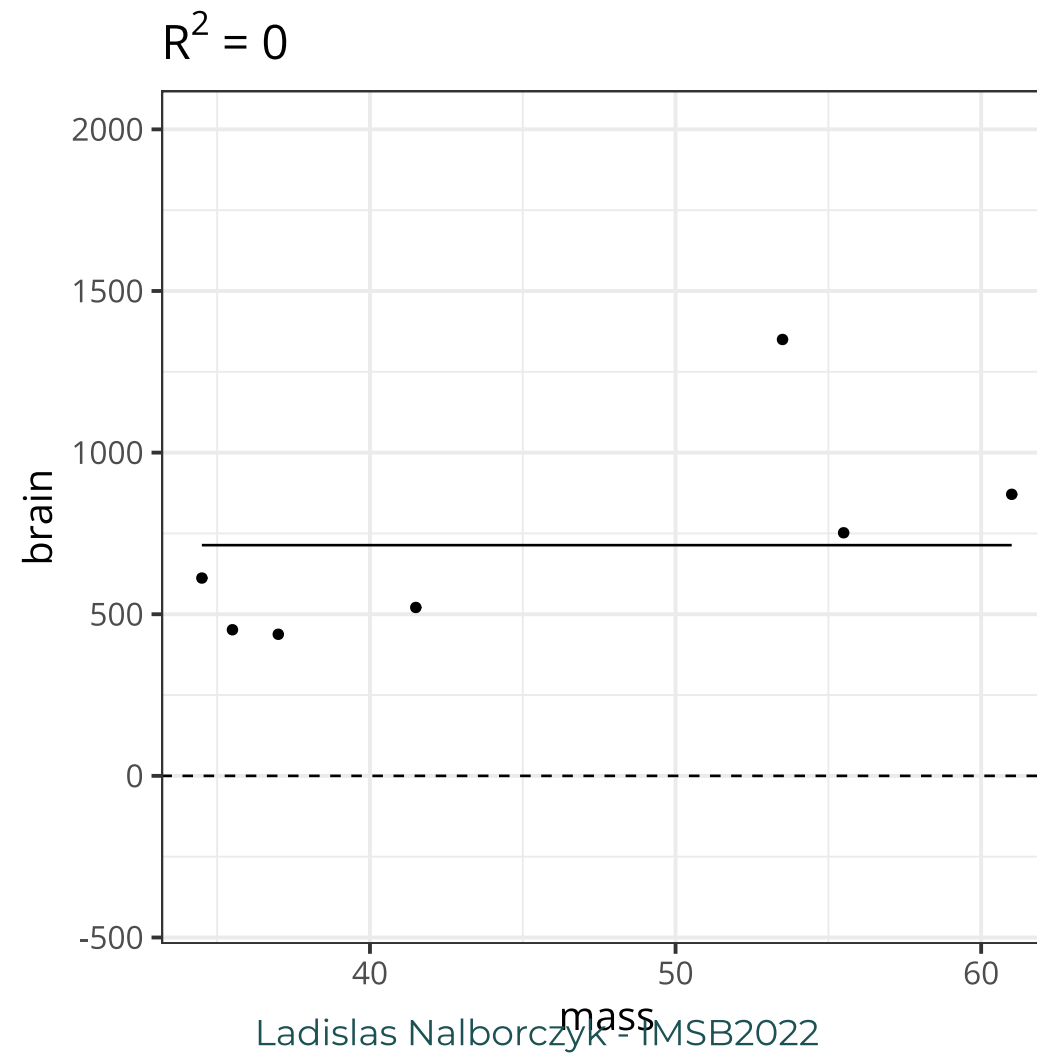


Underfitting

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

$$\mu_i = \alpha$$

```
1 mod1.7 <- lm(brain ~ 1, data = d)
```



Théorie de l'information

Il nous faut une autre mesure des capacités de prédiction pour évaluer nos modèles. Idéalement, on voudrait pouvoir mesurer la “distance” entre notre modèle et le “full model” (i.e., la “réalité”, la nature)... mais on ne sait pas faire cela.

Par contre, on peut mesurer de combien notre incertitude est réduite en découvrant un outcome (une observation) supplémentaire. Cette réduction est la définition de l'**information**.

Mais il nous faut tout d'abord une mesure de l'incertitude (pour savoir si on l'a réduite, ou pas)... S'il existe n événements possibles, et que chaque événement i a pour probabilité p_i , alors une mesure de l'incertitude est donnée par l'entropie (de Shannon) H :

$$H(p) = -\mathbb{E}[\log(p_i)] = \sum_{i=1}^n p_i \log\left(\frac{1}{p_i}\right) = -\sum_{i=1}^n p_i \log(p_i)$$

“

En d'autres termes : l'incertitude contenue dans une distribution de probabilités est la log-probabilité moyenne d'un événement.



Incertitude (exemple)

Exemple de prédiction météorologique. Si on imagine que la probabilité qu'il pleuve ou qu'il fasse beau (à Grenoble) est de $p_1 = 0.3$ et $p_2 = 0.7$.

Alors, $H(p) = -(p_1 \times \log(p_1) + p_2 \times \log(p_2)) \approx 0.61$.

```
1 p <- c(0.3, 0.7)
2 - sum(p * log(p) )
```

```
[1] 0.6108643
```

Imaginons que nous habitons à Abu Dhabi et que la probabilité qu'il y pleuve ou qu'il y fasse beau est de $p_1 = 0.01$ et $p_2 = 0.99$.

```
1 p <- c(0.01, 0.99)
2 - sum(p * log(p) )
```

```
[1] 0.05600153
```



Divergence

On a donc un moyen de quantifier l'incertitude. Comment utiliser cette mesure pour quantifier la distance entre notre modèle et la réalité ?

Divergence : incertitude ajoutée par l'utilisation d'une distribution de probabilités pour décrire... une autre distribution de probabilités ([Kullback-Leibler divergence](#), ou "entropie relative").

$$D_{\text{KL}}(p, q) = \sum_i p_i (\log(p_i) - \log(q_i)) = \sum_i p_i \log \left(\frac{p_i}{q_i} \right)$$

La divergence est la différence moyenne en log-probabilités entre la distribution cible (p) et le modèle (q).



Divergence

$$D_{KL}(p, q) = \sum_i p_i (\log(p_i) - \log(q_i)) = \sum_i p_i \log \left(\frac{p_i}{q_i} \right)$$

Par exemple, supposons que la “véritable” distribution de nos évènements (soleil vs pluie) soit $p_1 = 0.3$ et $p_2 = 0.7$. Si nous pensons plutôt que ces évènements arrivent avec une probabilité $q_1 = 0.25$ et $q_2 = 0.75$, quelle quantité d’incertitude avons-nous ajoutée ?

```
1 p <- c(0.3, 0.7)
2 q <- c(0.25, 0.75)
3
4 sum(p * log(p / q) )
```

```
[1] 0.006401457
```

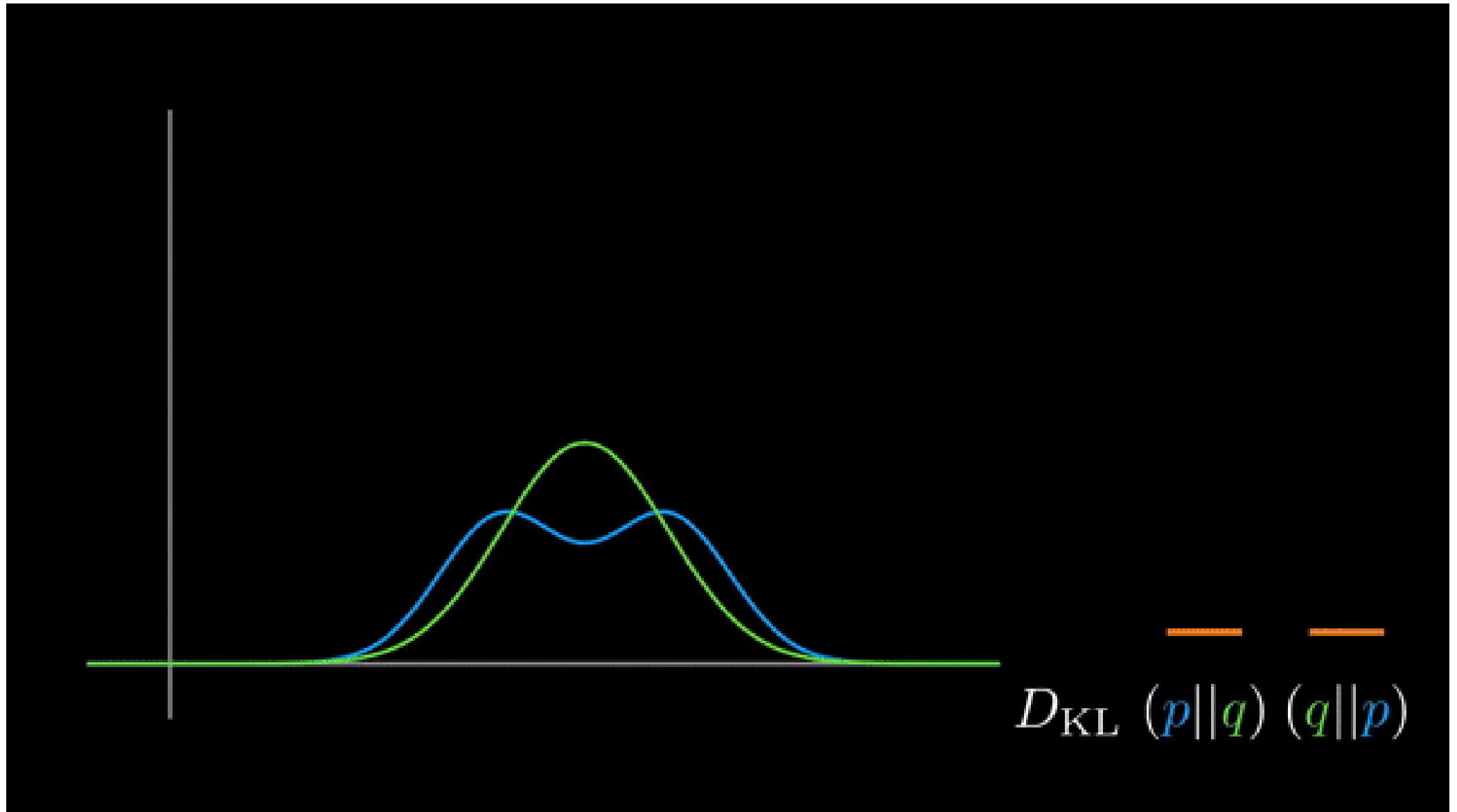
```
1 # NB : La divergence n'est pas symétrique...
2 sum(q * log(q / p) )
```

```
[1] 0.006164264
```



Divergence

La divergence n'est pas symétrique (ce n'est pas une distance)...



Entropie croisée et divergence

Entropie croisée : $H(p, q) = \sum_i p_i \log(q_i)$

```
1 sum(p * (log(q) ) )
```

La **Divergence** est définie comme l'entropie additionnelle ajoutée en utilisant q pour décrire p .

$$\begin{aligned} D_{\text{KL}}(p, q) &= H(p, q) - H(p) \\ &= - \sum_i p_i \log(q_i) - \left(- \sum_i p_i \log(p_i) \right) \\ &= - \sum_i p_i (\log(q_i) - \log(p_i)) \end{aligned}$$

```
1 - sum (p * (log(q) - log(p) ) )
```

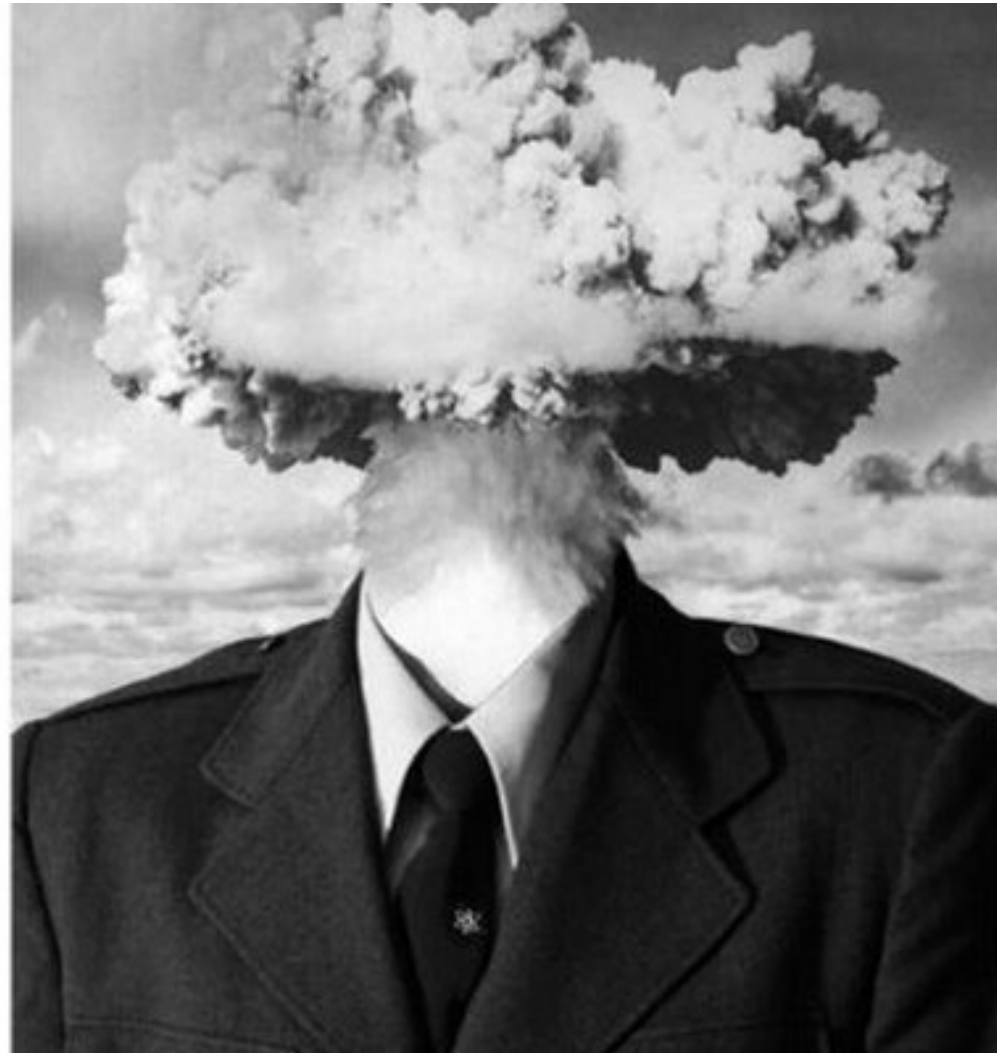
```
[1] 0.006401457
```



Vers la déviance...

OK, mais nous ne connaissons pas la distribution “cible” (la réalité), à quoi cela peut donc nous servir ?

Astuce : si nous comparons deux modèles, q et r , pour approximer p , nous allons comparer leurs divergences... Et donc $\mathbb{E}[\log(p_i)]$ sera la même quantité pour les deux modèles !



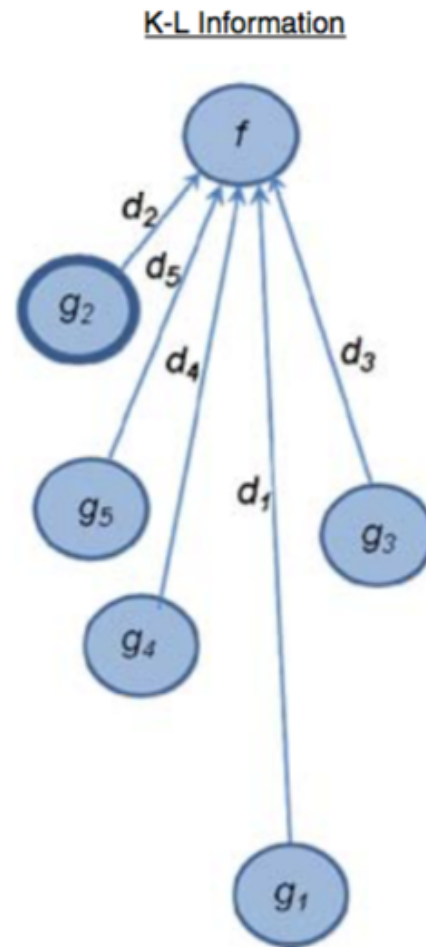
Ladislav Nalborczyk - IMSB2022





Vers la déviance...

On peut donc utiliser $\mathbb{E}[\log(q_i)]$ et $\mathbb{E}[\log(r_i)]$ comme estimateurs de la distance **relative** entre chaque modèle et notre distribution cible. On a seulement besoin de la **log-probabilité moyenne des modèles**. Comme on ne connaît pas la distribution cible, cela veut dire qu'on ne peut pas interpréter cette quantité en termes absolus mais seulement en termes relatifs. Ce qui nous intéresse c'est $\mathbb{E}[\log(q_i)] - \mathbb{E}[\log(r_i)]$.



Déviance

Pour approximer la valeur de $\mathbb{E}[\log(q_i)]$, on peut utiliser la déviance d'un modèle, qui est une mesure du fit **relatif** du modèle.

$$D(q) = -2 \sum_i \log(q_i)$$

où i indice chaque observation et q_i est la **vraisemblance** de chaque observation.

```
1 d$mass.s <- scale(d$mass)
2 mod1.8 <- lm(brain ~ mass.s, data = d)
3
4 -2 * logLik(mod1.8) # calcul de la déviance
```

```
'log Lik.' 94.92499 (df=3)
```



Déviance

```
1 # paramètres estimés (intercept et pente)
2 alpha <- coef(mod1.8)[1]
3 beta <- coef(mod1.8)[2]
4
5 # calcul de la log-vraisemblance
6 ll <- sum(dnorm(
7   x = d$brain,
8   mean = alpha + beta * d$mass.s,
9   sd = sigma(mod1.8),
10  log = TRUE
11  ) )
12
13 # calcul de la déviance
14 (-2) * ll
```

```
[1] 95.2803
```



Log-pointwise-predictive density

Les fréquentistes multiplient le log-score par -2 car la différence de deux déviations suit une loi de χ^2 , ce qui est utile pour tester l'hypothèse nulle. Mais sans besoin de tester l'hypothèse nulle (avec une forme prédéfinie), on peut très bien travailler directement avec le log-score $S(q) = \sum_i \log(q_i)$, qu'on traite comme une estimation de $\mathbb{E}[\log(q_i)]$.

On peut calculer $S(q)$ sur toute la distribution postérieure, ce qui donne la version bayésienne du log-score, la **log-pointwise-predictive density** :

$$\text{lppd}(y, \Theta) = \sum_i \log \frac{1}{S} \sum_s p(y_i \mid \Theta_s)$$

Où S est le nombre d'échantillons et Θ_s est le s -ième ensemble de valeurs de paramètres échantillonnés de la distribution postérieure.



In-sample and out-of-sample

La déviance a le même problème que le R^2 , lorsqu'elle est calculée sur l'échantillon observé. Dans ce cas, on l'appelle déviance **in-sample**.

Si on est intéressé par les capacités de prédiction de notre modèle, nous pouvons calculer la déviance du modèle sur de nouvelles données... qu'on appellera dans ce cas déviance **out-of-sample**. Cela revient à se demander si notre modèle est performant pour prédire de nouvelles données (non observées).

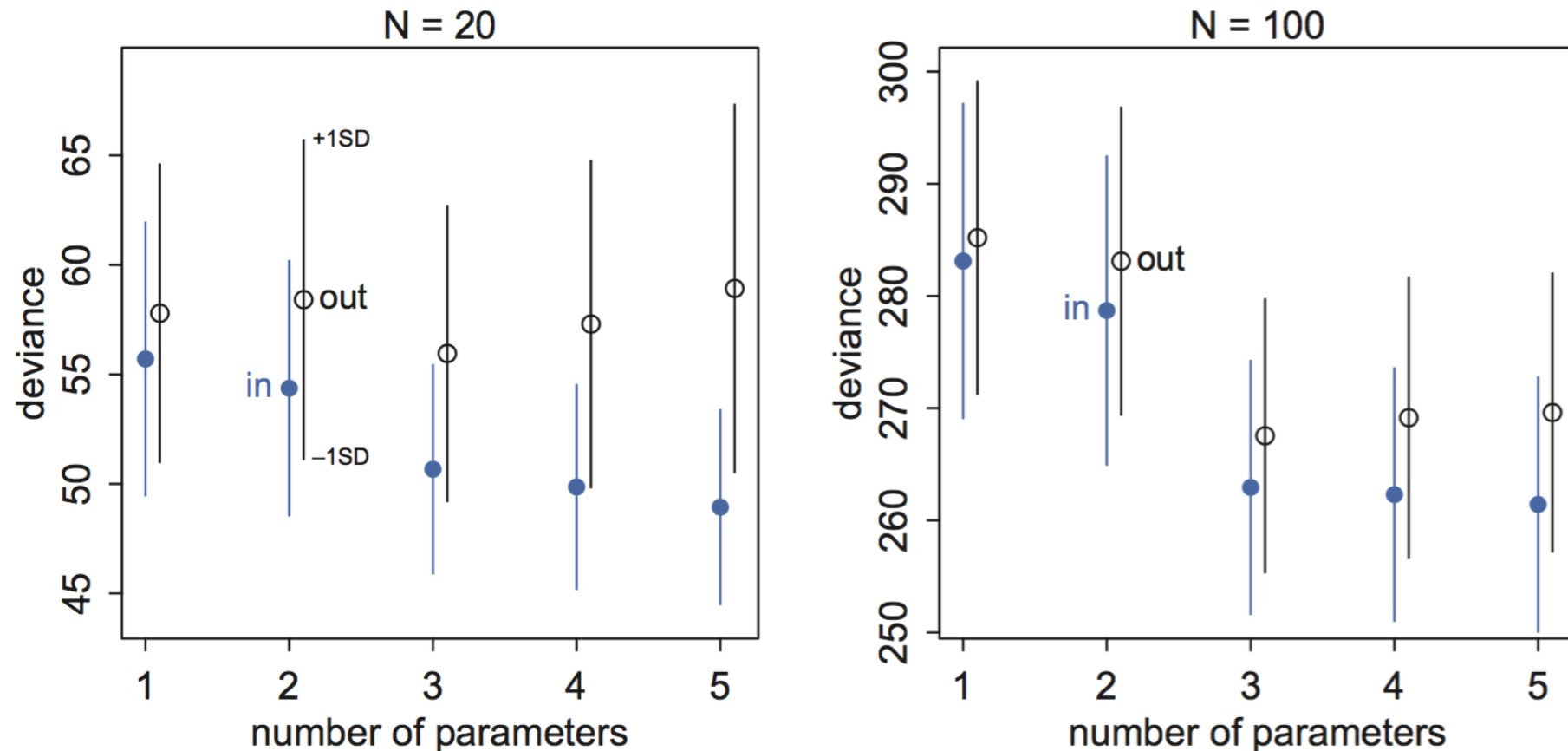
Imaginons que nous disposions d'un échantillon de taille N , que nous appellerons échantillon d'apprentissage (training). Nous pouvons calculer la déviance du modèle sur cet échantillon (D_{train} ou D_{in}). Si nous acquérons ensuite un nouvel échantillon de taille N issu du même processus de génération de données (que nous appellerons échantillon de test), nous pouvons calculer une déviance sur ce nouvel échantillon, en utilisant les paramètres estimés avec l'échantillon d'entraînement (que nous appellerons D_{test} ou D_{out}).



In sample and out of sample deviance

$$y_i \sim \text{Normal}(\mu_i, 1)$$

$$\mu_i = (0.15)x_{1,i} - (0.4)x_{2,i}$$

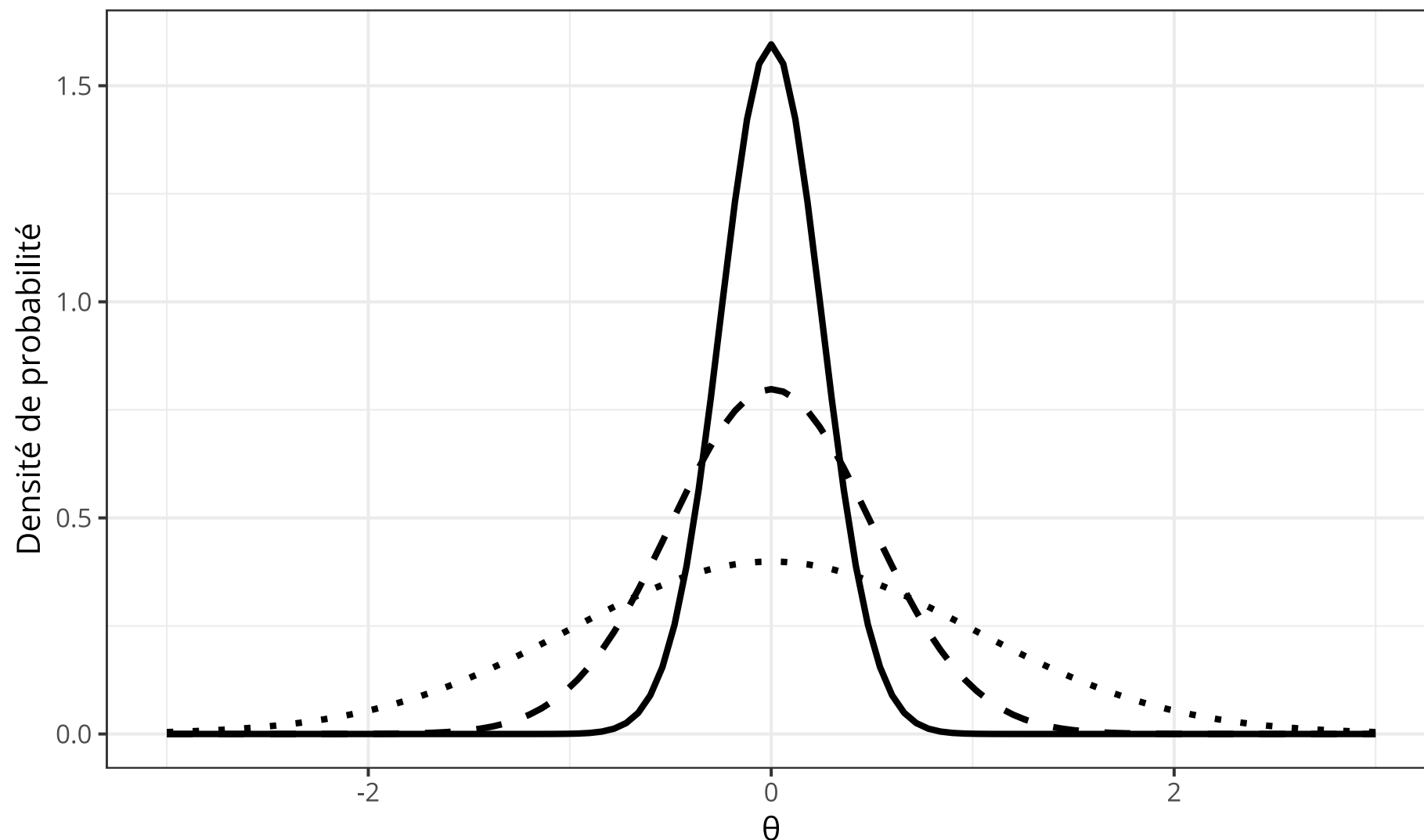


On a réalisé ce processus 10.000 fois pour cinq modèles de régression linéaire de complexité croissante. Les points bleus représentent la déviance calculée sur l'échantillon d'apprentissage et les points noirs la déviance calculée sur l'échantillon de test.

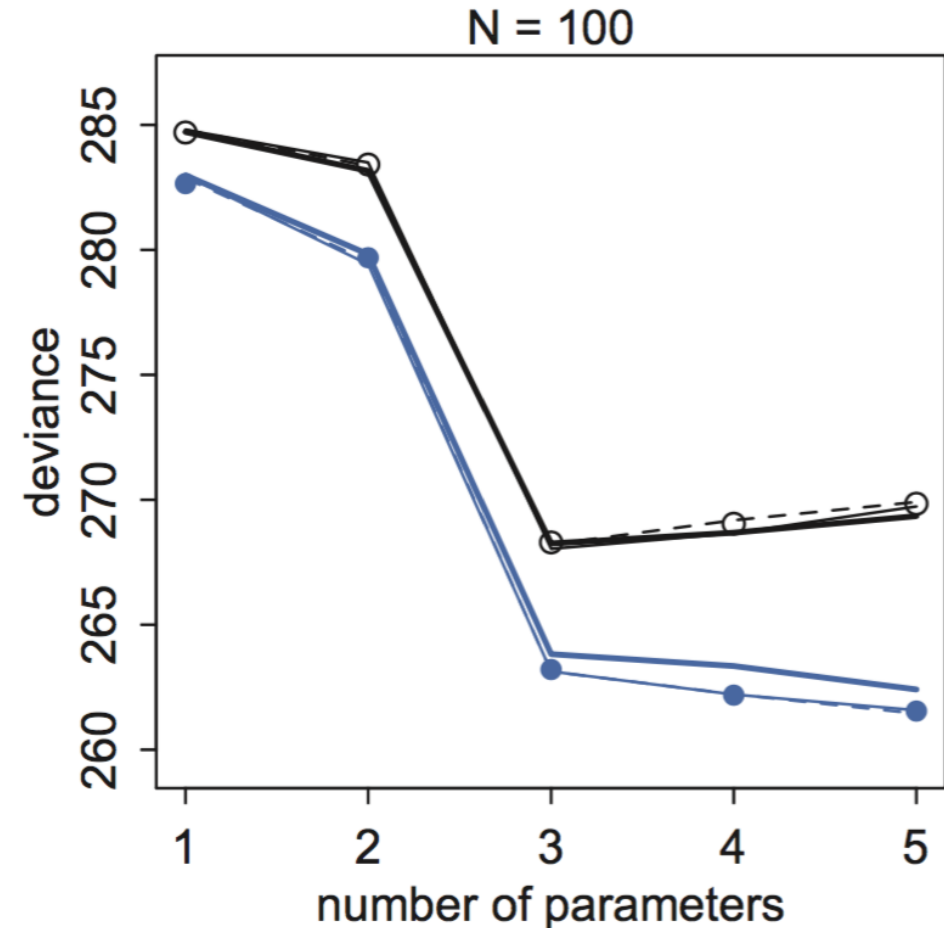
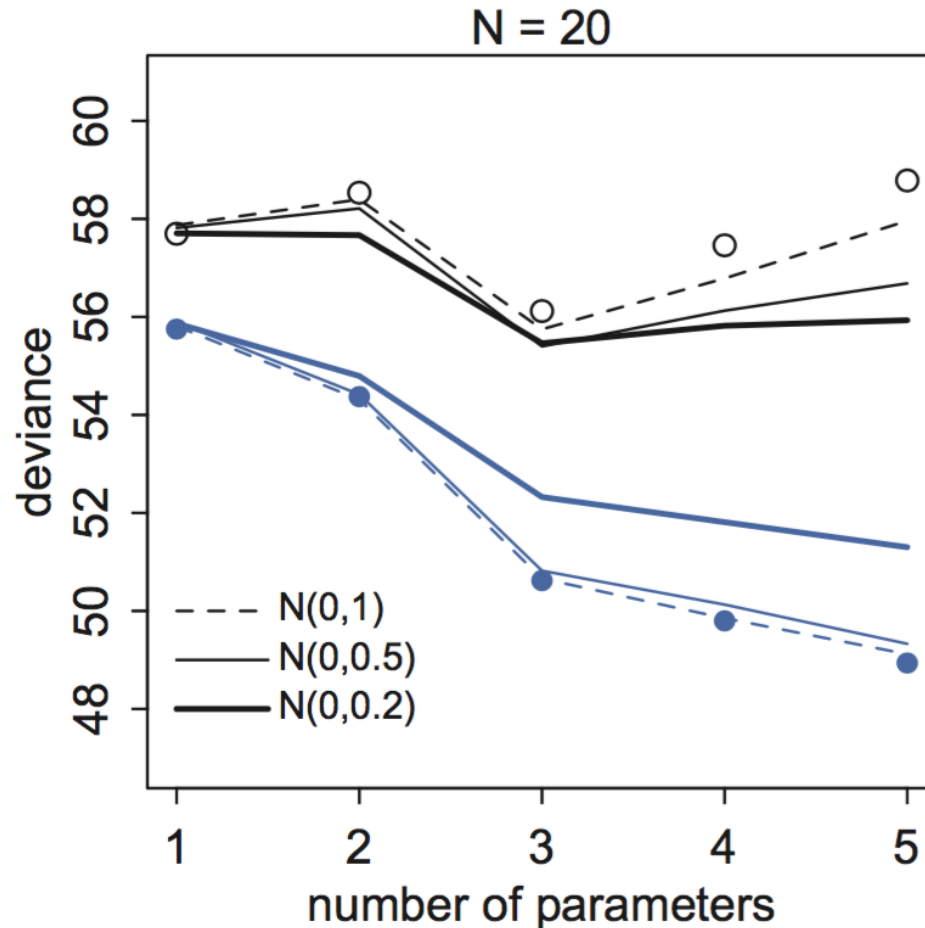


Régularisation

Une autre manière de lutter contre le sur-apprentissage (overfitting) est d'utiliser des priors “sceptiques” qui vont venir ralentir l'apprentissage réalisé sur les données (i.e., accorder plus de poids au prior).

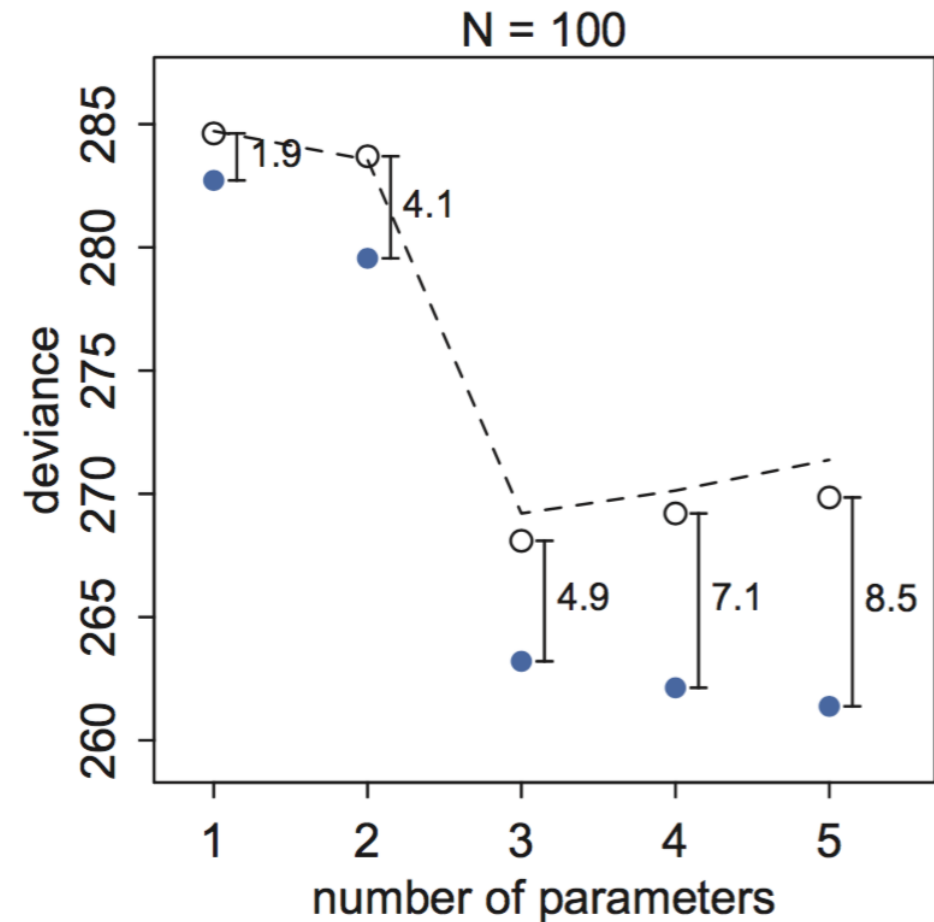
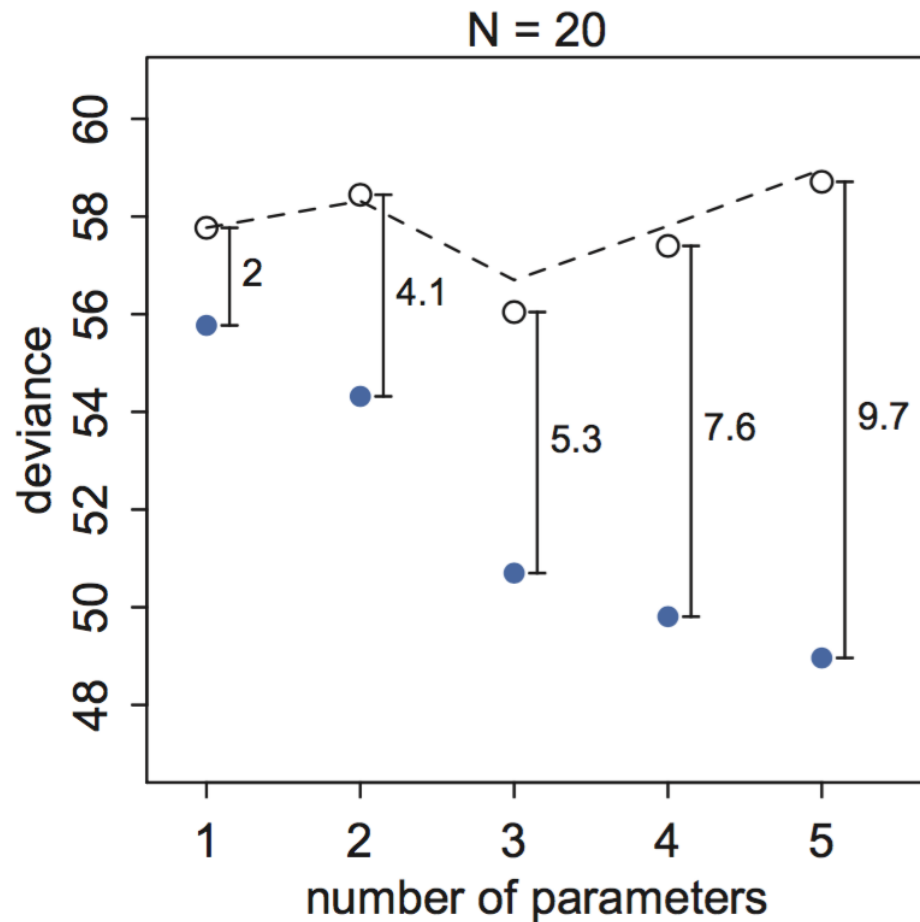


Régularisation



Comment décider de la précision du prior ? Est-ce que le prior est “assez” régularisateur ou pas ? On peut diviser le jeu de données en deux parties (training et testing) afin de choisir le prior qui produit la déviance **out-of-sample** la plus faible. On appelle cette stratégie la **validation croisée** (cross-validation).

Critères d'information



On mesure ici la différence entre la déviance **in-sample** (en bleu) et la déviance **out-of-sample** (en noir). On remarque que la déviance **out-of-sample** est presque exactement égale à la déviance **in-sample**, plus deux fois le nombre de paramètres du modèle...

Akaike information criterion

L'AIC fournit une approximation de la déviance **out-of-sample** :

$$AIC = D_{\text{train}} + 2p = -2\text{lppd} + 2p \approx D_{\text{test}}$$

où p est le nombre de paramètres libres (i.e., à estimer) dans le modèle. L'AIC donne donc une approximation des capacités de prédiction (out-of-sample) du modèle.¹



1. NB : l'AIC fonctionne bien uniquement quand le nombre d'observations N est largement supérieur au nombre de paramètres p . Dans le cas contraire, on utilise plutôt l'AICc ([Burnham & Anderson, 2004](#),
Ladislav Nalborczyk - IMSB2022



Deviance information criterion

Une condition d'application de l'AIC est que les priors soient plats ou dépassés par la vraisemblance (e.g., lorsqu'on a beaucoup de données). Le DIC est un indice qui ne requiert pas cette condition, en s'accommodant de priors informatifs.

Le DIC est calculé à partir de la distribution a posteriori de la déviance D calculée sur l'échantillon d'apprentissage (i.e., D_{train}).

$$\text{DIC} = \bar{D} + (\bar{D} - \hat{D}) = \bar{D} + p_D$$

où \bar{D} est la moyenne de la distribution a posteriori D calculée pour chaque valeur de paramètre échantillonnée, et \hat{D} la déviance calculée à la moyenne de la distribution a posteriori. La différence $\bar{D} - \hat{D} = p_D$ est analogue au nombre de paramètres utilisé dans le calcul de l'AIC (en cas de prior plat, cette différence revient à compter le nombre de paramètres).



Widely applicable information criterion

Une condition d'application de l'AIC et du DIC est que la distribution a posteriori soit une distribution gaussienne multivariée. Le WAIC relâche cette condition et il est souvent plus précis que le DIC.

Un aspect important du WAIC est qu'il est dit "pointwise", c'est à dire qu'il considère l'imprécision de prédiction point par point (donnée par donnée), indépendamment pour chaque observation.

On va commencer par calculer la **log-pointwise-predictive-density** (lppd), définie de la manière suivante :

$$\text{lppd}(y, \Theta) = \sum_i \log \frac{1}{S} \sum_s p(y_i | \Theta_s)$$

En Français : la log-densité prédictive point par point est **la somme du log de la vraisemblance moyenne de chaque observation**. Il s'agit de l'analogie point par point de la déviance, moyennée sur toute la distribution postérieure.



Widely applicable information criterion

```
1 library(brms)
2 data(cars)
3
4 priors <- c(
5   prior(normal(0, 100), class = Intercept),
6   prior(normal(0, 10), class = b),
7   prior(exponential(0.1), class = sigma)
8 )
9
10 mod1 <- brm(
11   dist ~ 1 + speed,
12   prior = priors,
13   data = cars
14 )
```



Widely applicable information criterion

$$\text{lppd}(y, \Theta) = \sum_i \log \frac{1}{S} \sum_s p(y_i | \Theta_s)$$

```

1 n_obs <- nrow(cars) # number of observations
2
3 ll <-
4   log_lik(mod1) %>% # pointwise log-likelihood (S rows * N observations)
5   data.frame() %>% # converts to dataframe
6   set_names(c(str_c(0, 1:9), 10:n_obs) ) # renaming columns
7
8 (lppd <-
9   ll %>%
10  pivot_longer(
11    everything(), # for all columns
12    names_to = "i",
13    values_to = "loglikelihood"
14  ) %>%
15  # log-likelihood to likelihood
16  mutate(likelihood = exp(loglikelihood) ) %>%
17  group_by(i) %>%
18  # taking the log of the average likelihood
19  summarise(log_mean_likelihood = mean(likelihood) %>% log() ) %>%
20  # computing the sum of these values
21  summarise(lppd = sum(log_mean_likelihood) ) %>%
22  pull(lppd) )

```

```
[1] -206.5503
```



Widely applicable information criterion

La deuxième partie du calcul du WAIC est le nombre de paramètres effectif, p_{WAIC} . On définit $\text{var}_\theta \log[p(y_i | \theta)]$ comme la variance de la log-vraisemblance pour chaque observation i de l'échantillon d'entraînement.

$$p_{WAIC} = \sum_i \text{var}_\theta \log[p(y_i | \theta)]$$

```

1 (pwaic <-
2   ll %>%
3   pivot_longer(
4     everything(),
5     names_to = "i",
6     values_to = "loglikelihood"
7   ) %>%
8   group_by(i) %>%
9   summarise(var_loglikelihood = var(loglikelihood) ) %>%
10  summarise(pwaic = sum(var_loglikelihood) ) %>%
11  pull() )

```

```
[1] 3.469262
```



Widely applicable information criterion

Ensuite, le WAIC est défini par :

$$\text{WAIC}(y, \Theta) = -2 \left(\text{lppd} - \underbrace{\sum_i \text{var}_{\theta} \log[p(y_i | \theta)]}_{\text{penalty term}} \right)$$

```
1 (WAIC <- -2 * (lppd - pwaic) )
```

```
[1] 420.0392
```



Widely applicable information criterion

Le WAIC est également un estimateur de la déviance **out-of-sample**. La fonction `brms::waic()` permet de le calculer directement.

```
1 waic(mod1)
```

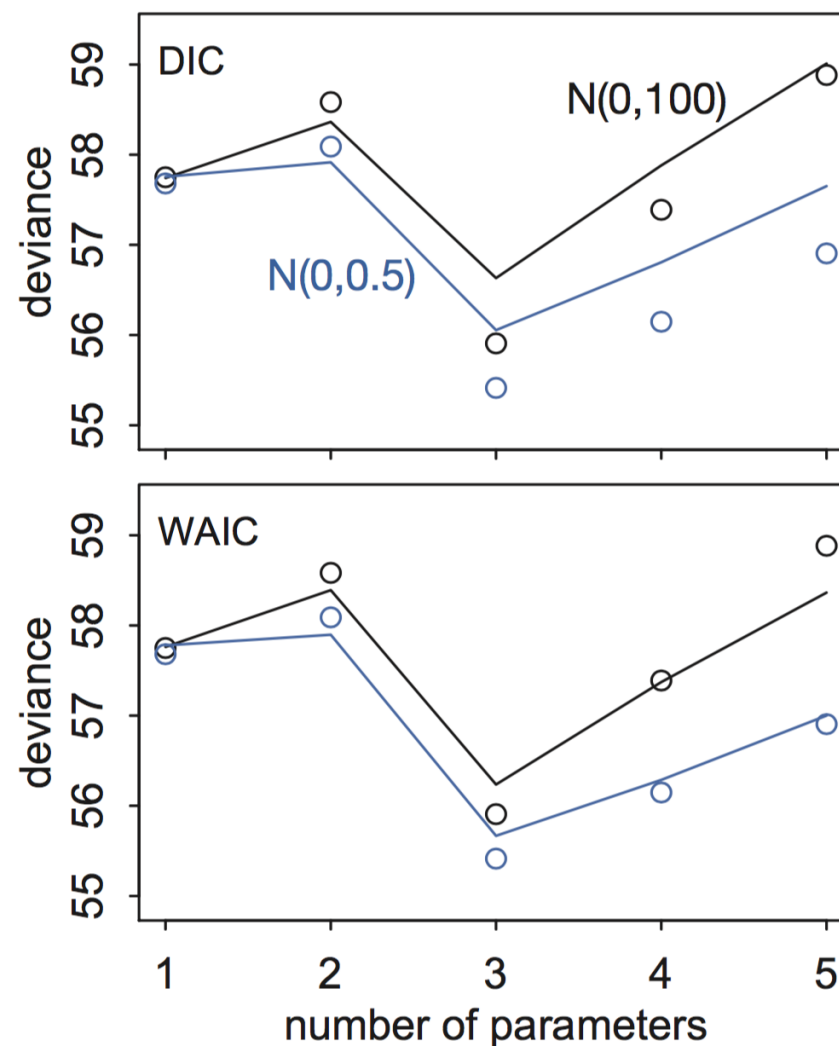
```
Computed from 4000 by 50 log-likelihood matrix
```

	Estimate	SE
elpd_waic	-210.0	6.6
p_waic	3.5	1.3
waic	420.0	13.2

```
2 (4.0%) p_waic estimates greater than 0.4. We recommend trying loo instead.
```



Critères d'information et régularisation



Le DIC et le WAIC peuvent être conceptualisés (au même titre que l'AIC) comme des approximations de la déviance **out-of-sample**. On remarque que le WAIC produit des approximations plus précises que le DIC, et que l'utilisation de priors régularisateurs permet de réduire la déviance **out-of-sample**.



Et ensuite ?

Sélection de modèle : On choisit le meilleur modèle en utilisant un des outils présentés et on base nos conclusions sur les paramètres estimés par ce meilleur modèle.

Comparaison de modèles : On utilise la validation croisée ou des critères d'informations mais aussi les outils de **posterior predictive checking** discutés précédemment, pour chaque modèle, afin d'étudier leurs forces et faiblesses.

Moyennage de modèles : On va construire des posterior predictive checks qui exploitent ce qu'on sait des capacités de prédiction de chaque modèle (e.g., via le WAIC).



Comparaison de modèles

On essaye de prédire les kg par gramme de lait (**kcal.per.g**) avec les prédicteurs **neocortex** et le logarithme de **mass**. Nous allons ensuite fitter 4 modèles qui correspondent aux 4 combinaisons possibles de prédicteurs et les comparer en utilisant le WAIC.

```
1 library(imsb)
2 d <- open_data(milk)
3 d <- milk[complete.cases(milk), ] # removing NAs
4 d$neocortex <- d$neocortex.perc / 100 # rescaling explanatory variable
5 head(d)
```

	clade	species	kcal.per.g	perc.fat	perc.protein
1	Strepsirrhine	Eulemur fulvus	0.49	16.60	15.42
6	New World Monkey	Alouatta seniculus	0.47	21.22	23.58
7	New World Monkey	A palliata	0.56	29.66	23.46
8	New World Monkey	Cebus apella	0.89	53.41	15.80
10	New World Monkey	S sciureus	0.92	50.58	22.33
11	New World Monkey	Cebuella pygmaea	0.80	41.35	20.85

	perc.lactose	mass	neocortex.perc	neocortex
1	67.98	1.95	55.16	0.5516
6	55.20	5.25	64.54	0.6454
7	46.88	5.37	64.54	0.6454
8	30.79	2.51	67.64	0.6764
10	27.09	0.68	68.85	0.6885
11	37.80	0.12	58.85	0.5885



Comparaison de modèles

```
1 mod2.1 <- brm(  
2   kcal.per.g ~ 1,  
3   family = gaussian,  
4   data = d,  
5   prior = c(  
6     prior(normal(0, 100), class = Intercept),  
7     prior(exponential(0.01), class = sigma)  
8   ),  
9   iter = 2000, warmup = 1000,  
10  backend = "cmdstanr" # on peut changer le backend de brms  
11 )  
12  
13 mod2.2 <- brm(  
14   kcal.per.g ~ 1 + neocortex,  
15   family = gaussian,  
16   data = d,  
17   prior = c(  
18     prior(normal(0, 100), class = Intercept),  
19     prior(normal(0, 10), class = b),  
20     prior(exponential(0.01), class = sigma)  
21   ),  
22   iter = 2000, warmup = 1000,  
23   backend = "cmdstanr"  
24 )
```



Comparaison de modèles

```
1 mod2.3 <- brm(  
2   kcal.per.g ~ 1 + log(mass),  
3   family = gaussian,  
4   data = d,  
5   prior = c(  
6     prior(normal(0, 100), class = Intercept),  
7     prior(exponential(0.01), class = sigma)  
8   ),  
9   iter = 2000, warmup = 1000,  
10  backend = "cmdstanr"  
11 )  
12  
13 mod2.4 <- brm(  
14   kcal.per.g ~ 1 + neocortex + log(mass),  
15   family = gaussian,  
16   data = d,  
17   prior = c(  
18     prior(normal(0, 100), class = Intercept),  
19     prior(normal(0, 10), class = b),  
20     prior(exponential(0.01), class = sigma)  
21   ),  
22   iter = 2000, warmup = 1000,  
23   backend = "cmdstanr"  
24 )
```



Comparaison de modèles

On peut utiliser la méthode `update()` qui permet de fitter plus rapidement un nouveau modèle qui ressemble à un modèle déjà existant.

```
1 mod2.3 <- update(  
2   object = mod2.2,  
3   newdata = d,  
4   formula = kcal.per.g ~ 1 + log(mass)  
5 )  
6  
7 mod2.4 <- update(  
8   object = mod2.3,  
9   newdata = d,  
10  formula = kcal.per.g ~ 1 + neocortex + log(mass)  
11 )
```



Comparaison de modèles

```
1 # calcul du WAIC et ajout du WAIC à chaque modèle
2
3 mod2.1 <- add_criterion(mod2.1, "waic")
4 mod2.2 <- add_criterion(mod2.2, "waic")
5 mod2.3 <- add_criterion(mod2.3, "waic")
6 mod2.4 <- add_criterion(mod2.4, "waic")
7
8 # comparaison des WAIC de chaque modèle
9
10 w <- loo_compare(mod2.1, mod2.2, mod2.3, mod2.4, criterion = "waic")
11 print(w, simplify = FALSE)
```

	elpd_diff	se_diff	elpd_waic	se_elpd_waic	p_waic	se_p_waic	waic	se_waic
mod2.4	0.0	0.0	8.4	2.6	3.1	0.8	-16.8	5.1
mod2.3	-3.9	1.7	4.5	2.1	2.0	0.4	-9.0	4.2
mod2.1	-4.0	2.4	4.4	1.9	1.3	0.3	-8.8	3.7
mod2.2	-4.8	2.5	3.6	1.6	1.9	0.3	-7.2	3.3



Akaike's weights

Le poids d'un modèle est une estimation de la probabilité que ce modèle fera les meilleures prédictions possibles sur un nouveau jeu de données, conditionnellement au set de modèles considéré.

$$w_i = \frac{\exp\left(-\frac{1}{2}d\text{WAIC}_i\right)}{\sum_{j=1}^m \exp\left(-\frac{1}{2}d\text{WAIC}_j\right)}$$

Cette fonction permet simplement de passer du WAIC à une probabilité (équivalent à une fonction softmax). Le modèle `mod2.4` a un poids de 0.96, qui le place en tête du jeu de modèles. N'oublions cependant pas que nous disposons seulement de 17 observations...

```
1 model_weights(mod2.1, mod2.2, mod2.3, mod2.4, weights = "waic") %>% round(digits = 3)
```

```
mod2.1 mod2.2 mod2.3 mod2.4
0.017  0.008  0.019  0.955
```



Moyennage de modèles (model averaging)

Pourquoi ne conserver uniquement le premier modèle et oublier les autres ? Une autre stratégie consisterait à pondérer les prédictions des modèles par leurs poids respectifs. C'est ce qu'on appelle le moyennage de modèles (model averaging).

- Calculer le WAIC de chaque modèle.
- Calculer le poids de chaque modèle.
- Simuler des données à partir de chaque modèle.
- Combiner ces valeurs simulées dans un **ensemble** de prédictions pondérées par le poids du modèle.



Moyennage de modèles (model averaging)

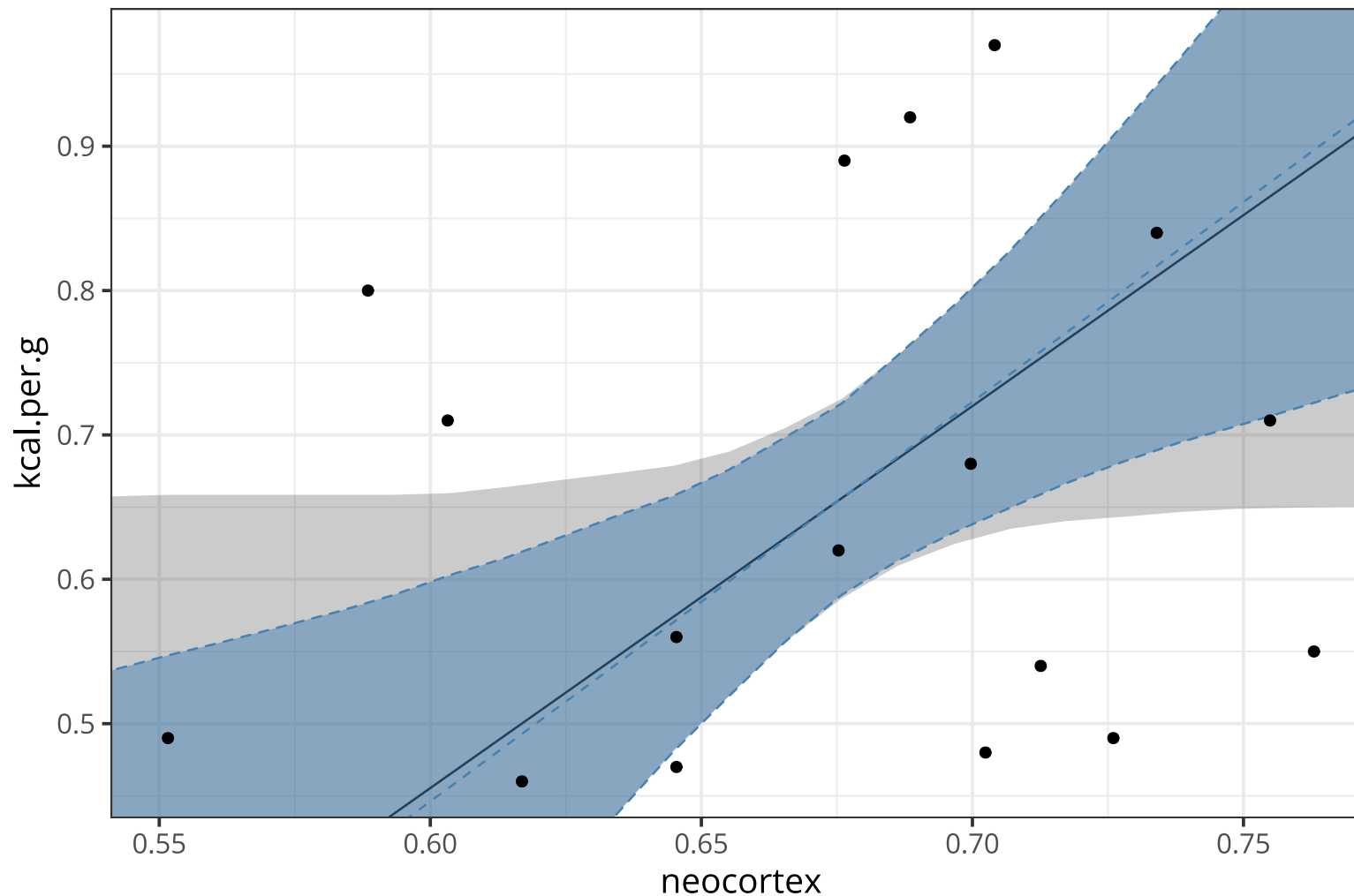
On peut utiliser la fonction `brms::pp_average()` qui pondère les prédictions de chaque modèle par leur poids.

```
1 # grille de valeurs pour lesquelles on va générer des prédictions
2 new_data <- data.frame(
3   neocortex = seq(from = 0.5, to = 0.8, length.out = 30),
4   mass = 4.5
5 )
6
7 # prédictions du modèle mod2.4
8 f <- fitted(mod2.4, newdata = new_data) %>%
9   as.data.frame() %>%
10  bind_cols(new_data)
11
12 # prédictions moyennées sur les 4 modèles
13 averaged_predictions <- pp_average(
14   mod2.1, mod2.2, mod2.3, mod2.4,
15   weights = "waic",
16   method = "fitted",
17   newdata = new_data
18 ) %>%
19   as.data.frame() %>%
20   bind_cols(new_data)
```



Moyennage de modèles (model averaging)

Voici les prédictions de tous les modèles considérés, pondérés par leur poids respectif. Comme le modèle `mod2.4` concentrait quasiment tout le poids, il fait sens que cette prédiction moyennée soit similaire aux prédictions du modèle `mod2.4`.

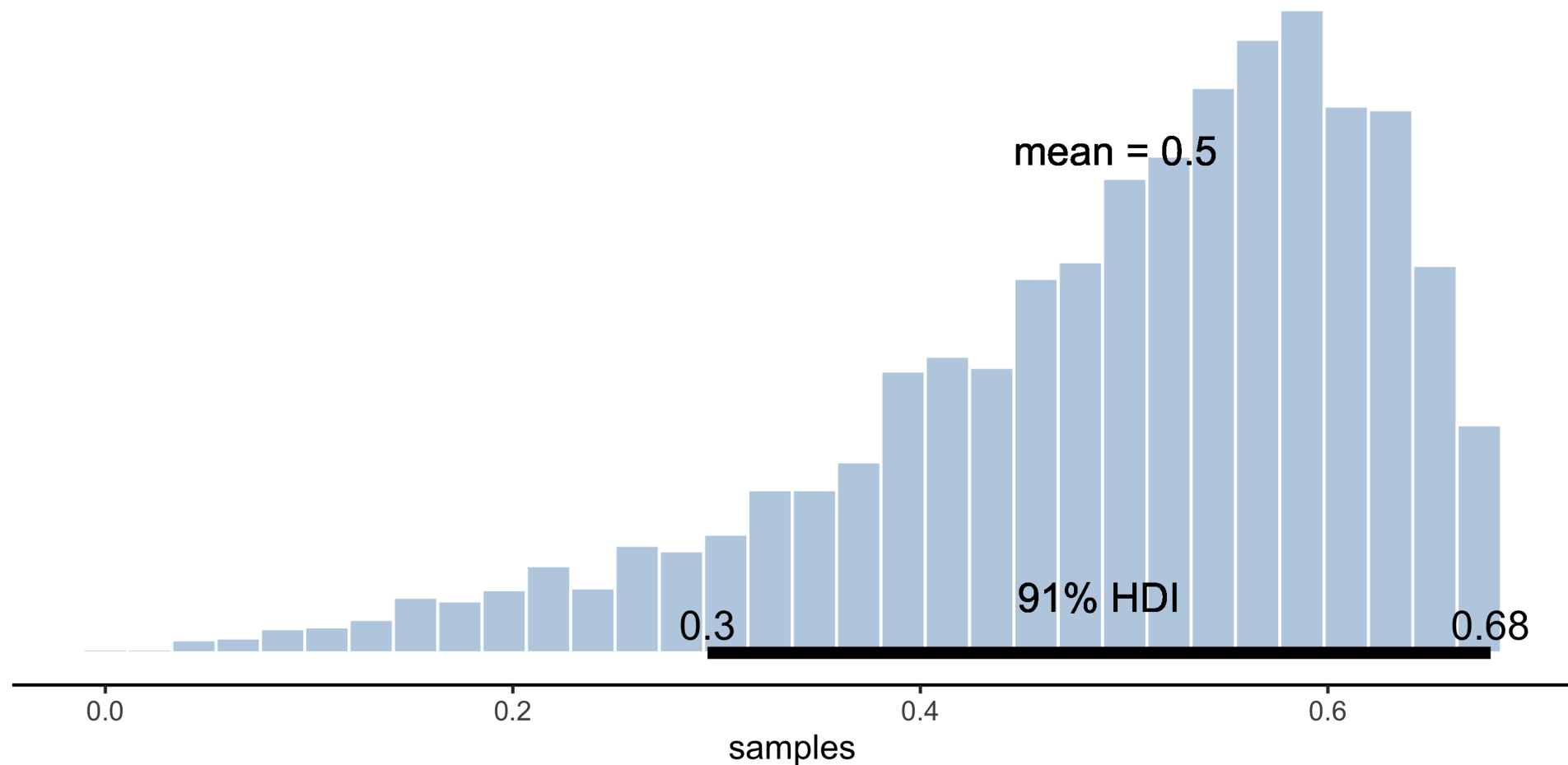


R-squared

```
1 bayes_R2(mod2.4) %>% round(digits = 3)
```

	Estimate	Est.Error	Q2.5	Q97.5
R2	0.496	0.13	0.167	0.667

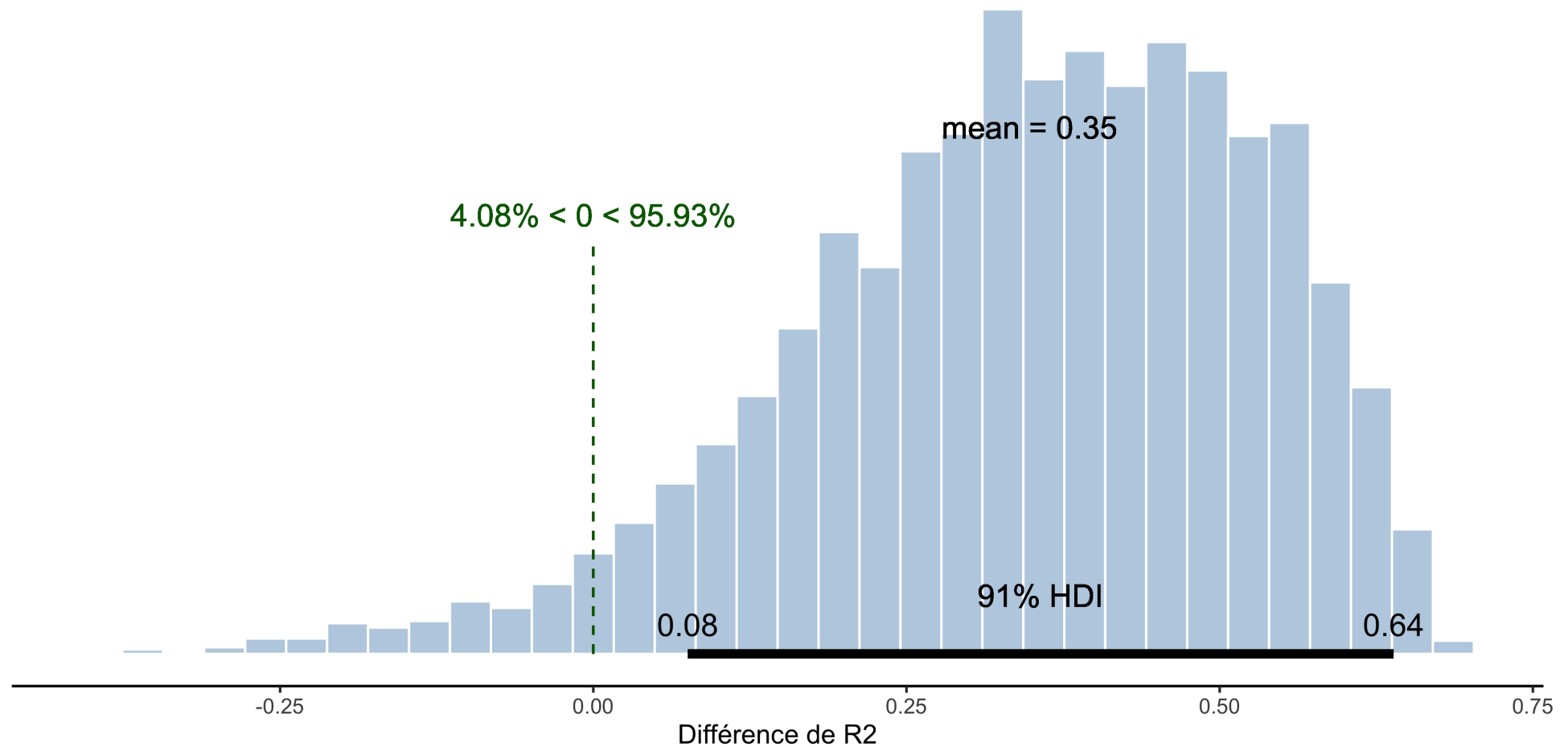
```
1 posterior_plot(samples = bayes_R2(mod2.4, summary = FALSE)[, 1])
```





R-squared

```
1 posterior_plot(  
2   samples = bayes_R2(mod2.4, summary = FALSE)[, 1] -  
3     bayes_R2(mod2.3, summary = FALSE)[, 1],  
4   compval = 0  
5 ) + labs(x = "Différence de R2")
```



Conclusions

Se méfier des interprétations intuitives (et souvent erronées) de la p-valeur, des intervalles de confiance, ou du facteur de Bayes.

Retenir que la difficulté en modélisation est de trouver un juste équilibre entre sous-apprentissage (underfitting) et sur-apprentissage (overfitting).

Pour contraindre l'apprentissage réalisé par le modèle sur les données observées (et ainsi éviter que le modèle accorde trop de poids à ces données), on peut utiliser des priors dits **régularisateurs** et/ou des outils comme la validation croisée ou les critères d'informations permettant d'estimer les capacités de prédiction du modèle sur de nouvelles données.



Travaux pratiques

```
1 # import des données "howell"
2 d <- open_data(howell) %>% mutate(age = scale(age) )
3
4 # on définit une graine (afin de pouvoir reproduire les résultats)
5 set.seed(666)
6
7 # on échantillonne des lignes du jeu de données
8 i <- sample(1:nrow(d), size = nrow(d) / 2)
9
10 # on définit l'échantillon d'entraînement
11 d1 <- d[i, ]
12
13 # on définit l'échantillon de test
14 d2 <- d[-i, ]
```

Nous avons maintenant deux dataframes, de 272 lignes chacune. On va utiliser **d1** pour fitter nos modèles et **d2** pour les évaluer.



Travaux pratiques

Soit h_i les valeurs de taille et x_i les valeurs centrées d'âge, sur la ligne i . Construisez les modèles suivants avec `d1`, en utilisant `brms::brm()` et des priors faiblement régularisateurs.

$$\mathcal{M}_1 : h_i \sim \text{Normal}(\mu_i, \sigma)$$

$$\mu_i = \alpha + \beta_1 x_i$$

$$\mathcal{M}_2 : h_i \sim \text{Normal}(\mu_i, \sigma)$$

$$\mu_i = \alpha + \beta_1 x_i + \beta_2 x_i^2$$

$$\mathcal{M}_3 : h_i \sim \text{Normal}(\mu_i, \sigma)$$

$$\mu_i = \alpha + \beta_1 x_i + \beta_2 x_i^2 + \beta_3 x_i^3$$

$$\mathcal{M}_4 : h_i \sim \text{Normal}(\mu_i, \sigma)$$

$$\mu_i = \alpha + \beta_1 x_i + \beta_2 x_i^2 + \beta_3 x_i^3 + \beta_4 x_i^4$$

$$\mathcal{M}_5 : h_i \sim \text{Normal}(\mu_i, \sigma)$$

$$\mu_i = \alpha + \beta_1 x_i + \beta_2 x_i^2 + \beta_3 x_i^3 + \beta_4 x_i^4 + \beta_5 x_i^5$$

$$\mathcal{M}_6 : h_i \sim \text{Normal}(\mu_i, \sigma)$$

$$\mu_i = \alpha + \beta_1 x_i + \beta_2 x_i^2 + \beta_3 x_i^3 + \beta_4 x_i^4 + \beta_5 x_i^5 + \beta_6 x_i^6$$



Travaux pratiques

- Comparer ces modèles en utilisant le WAIC. Comparer les rangs des modèles et leurs poids.
- Pour chaque modèle, produire un plot de la moyenne estimée et l'intervalle de confiance à 97% de la moyenne, surimposée aux données brutes. Comment ces prédictions diffèrent-elles selon les modèles ?
- Faire un plot des prédictions moyennées sur tous les modèles (sur les trois meilleurs). En quoi ces prédictions diffèrent-elles des prédictions du modèle avec le plus petit WAIC ?
- Calculer la déviance **out-of-sample** pour chaque modèle. Comparer les déviations obtenues à la question précédente aux valeurs de WAIC. Basé sur les déviations obtenues, quel modèle fait les meilleures prédictions ? Est-ce que le WAIC est un bon estimateur de la déviance **out-of-sample** ?



Solution - Question 1

```
1 mod3.1 <- brm(  
2   height ~ 1 + age,  
3   family = gaussian(),  
4   data = d1,  
5   prior = c(  
6     prior(normal(0, 100), class = Intercept),  
7     prior(exponential(0.01), class = sigma)  
8   ),  
9   backend = "cmdstanr"  
10  )  
11  
12 mod3.2 <- update(  
13   mod3.1,  
14   newdata = d1,  
15   formula = height ~ 1 + age + I(age^2)  
16  )  
17  
18 mod3.3 <- update(  
19   mod3.1,  
20   newdata = d1,  
21   formula = height ~ 1 + age + I(age^2) + I(age^3)  
22  )
```



Solution - Question 1

```
1 mod3.4 <- update(  
2   mod3.1,  
3   newdata = d1,  
4   formula = height ~ 1 + age + I(age^2) + I(age^3) + I(age^4)  
5 )  
6  
7 mod3.5 <- update(  
8   mod3.1,  
9   newdata = d1,  
10  formula = height ~ 1 + age + I(age^2) + I(age^3) + I(age^4) + I(age^5)  
11 )  
12  
13 mod3.6 <- update(  
14   mod3.1,  
15   newdata = d1,  
16   formula = height ~ 1 + age + I(age^2) + I(age^3) + I(age^4) + I(age^5) + I(age^6)  
17 )
```



Solution - Question 1

```

1 # calcul du WAIC et ajout du WAIC à chaque modèle
2
3 mod3.1 <- add_criterion(mod3.1, "waic")
4 mod3.2 <- add_criterion(mod3.2, "waic")
5 mod3.3 <- add_criterion(mod3.3, "waic")
6 mod3.4 <- add_criterion(mod3.4, "waic")
7 mod3.5 <- add_criterion(mod3.5, "waic")
8 mod3.6 <- add_criterion(mod3.6, "waic")
9
10 # comparaison des WAIC de chaque modèle
11
12 mod_comp <- loo_compare(mod3.1, mod3.2, mod3.3, mod3.4, mod3.5, mod3.6, criterion = "waic")
13 print(mod_comp, digits = 2, simplify = FALSE)

```

	elpd_diff	se_diff	elpd_waic	se_elpd_waic	p_waic	se_p_waic	waic
mod3.4	0.00	0.00	-954.40	12.41	5.74	0.79	1908.79
mod3.5	-0.70	0.47	-955.10	12.48	6.46	0.87	1910.20
mod3.6	-1.29	1.03	-955.69	12.27	7.44	0.93	1911.37
mod3.3	-21.90	6.37	-976.30	11.83	6.11	1.31	1952.60
mod3.2	-132.62	13.63	-1087.02	11.29	5.51	1.28	2174.03
mod3.1	-259.20	15.32	-1213.60	10.91	3.43	0.42	2427.20

	se_waic
mod3.4	24.83
mod3.5	24.95
mod3.6	24.53
mod3.3	23.66
mod3.2	22.58
mod3.1	21.82



Solution - Question 2

```

1 # on crée un vecteur de valeurs possibles pour "age"
2 age_seq <- data.frame(age = seq(from = -2, to = 3, length.out = 1e2) )
3
4 # on récupère les prédictions du modèle pour ces valeurs
5 mu <- data.frame(fitted(mod3.1, newdata = age_seq) ) %>% bind_cols(age_seq)
6
7 # on récupère les prédictions du modèle pour ces valeurs
8 pred_age <- data.frame(predict(mod3.1, newdata = age_seq) ) %>% bind_cols(age_seq)
9
10 # on affiche les dix premières prédictions
11 head(pred_age, 10)

```

	Estimate	Est.Error	Q2.5	Q97.5	age
1	100.4499	21.13775	58.51867	142.4313	-2.000000
2	101.0368	20.87114	60.20049	141.1053	-1.949495
3	101.8842	20.93198	60.09004	141.7364	-1.898990
4	103.0255	21.18667	60.41278	144.7774	-1.848485
5	104.5247	21.22701	63.20084	145.5188	-1.797980
6	104.9453	20.88182	65.70421	146.6455	-1.747475
7	105.6927	20.74311	64.78789	145.7145	-1.696970
8	106.8630	20.86651	65.64000	148.4519	-1.646465
9	107.6850	20.79897	66.01354	148.0521	-1.595960
10	108.9927	20.77924	67.99611	149.7180	-1.545455

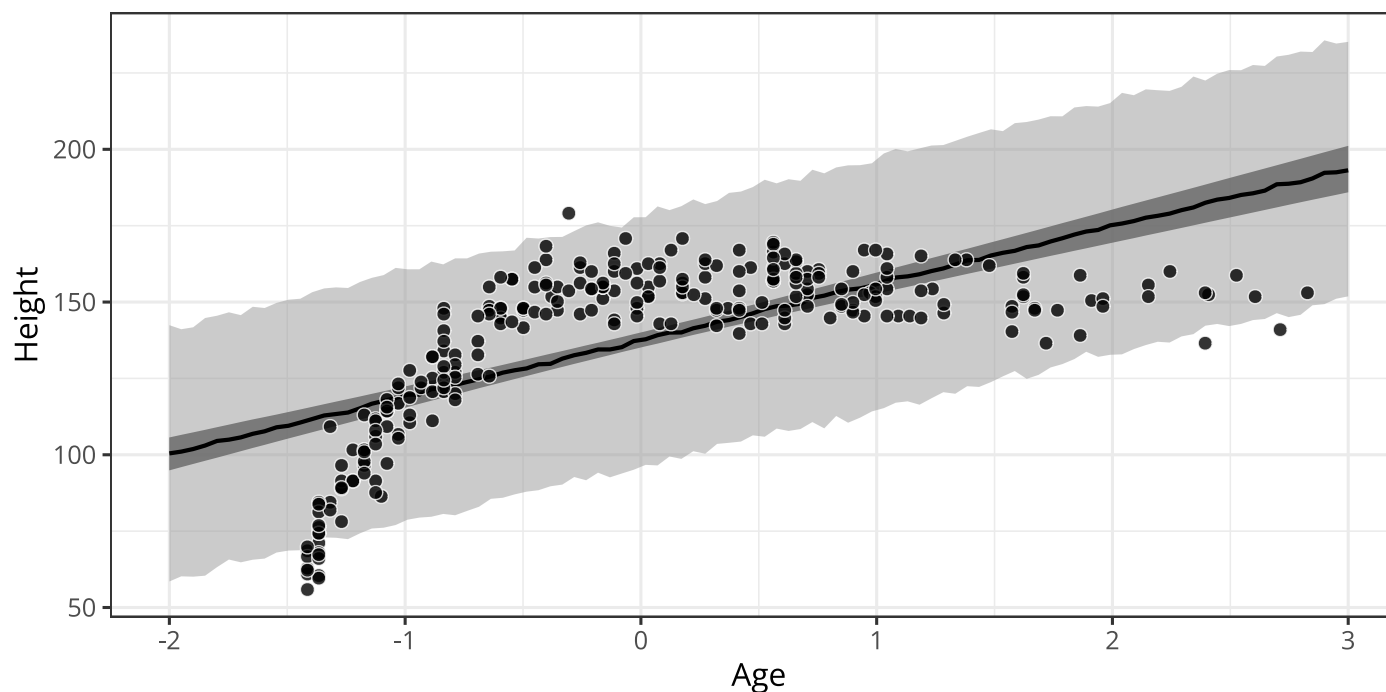


Solution - Question 2

```

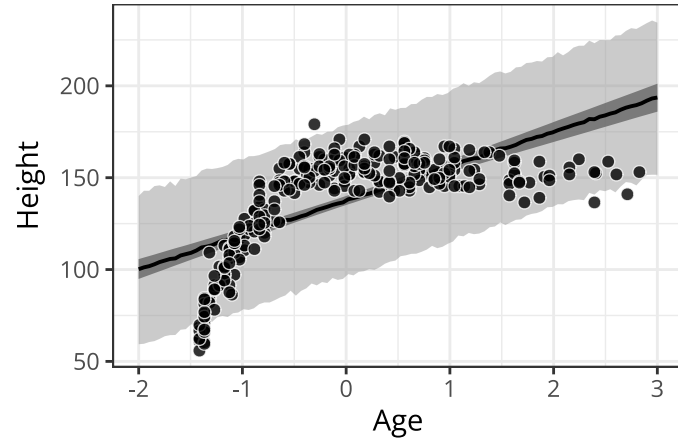
1 d1 %>%
2   ggplot(aes(x = age, y = height) ) +
3   geom_ribbon(
4     data = mu, aes(x = age, ymin = Q2.5, ymax = Q97.5),
5     alpha = 0.8, inherit.aes = FALSE
6   ) +
7   geom_smooth(
8     data = pred_age, aes(y = Estimate, ymin = Q2.5, ymax = Q97.5),
9     stat = "identity", color = "black", alpha = 0.5, size = 1
10  ) +
11  geom_point(colour = "white", fill = "black", pch = 21, size = 3, alpha = 0.8) +
12  labs(x = "Age", y = "Height")

```

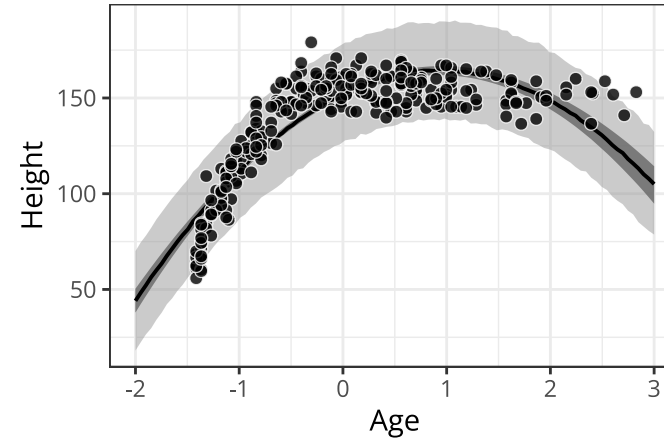


Solution - Question 2

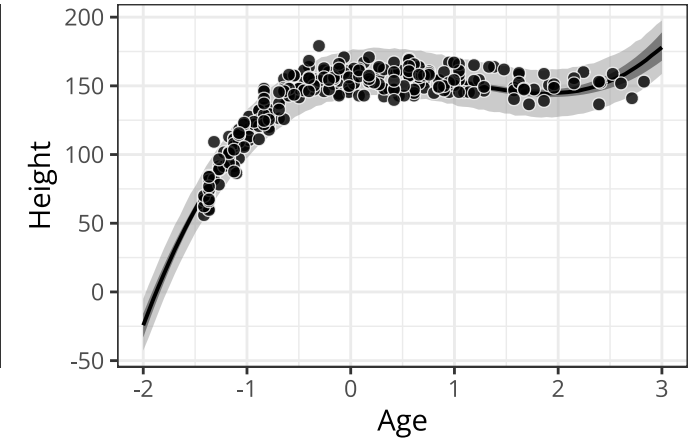
Predictions of mod3.1



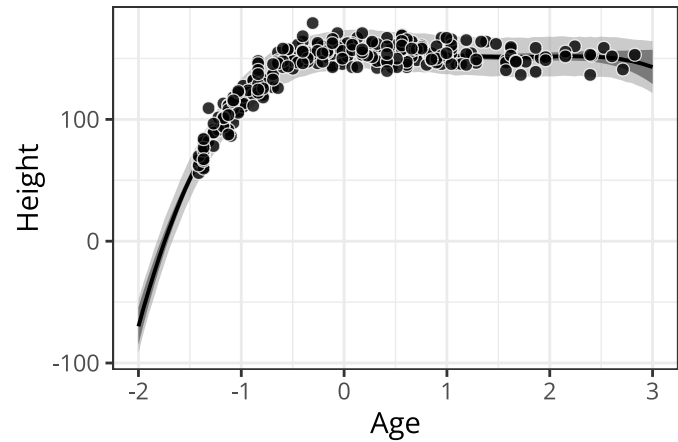
Predictions of mod3.2



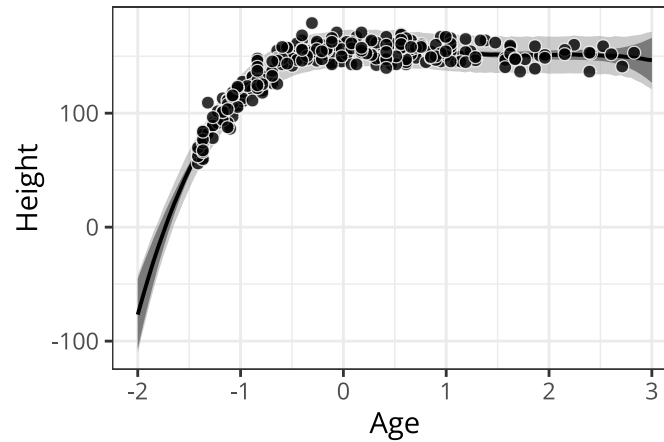
Predictions of mod3.3



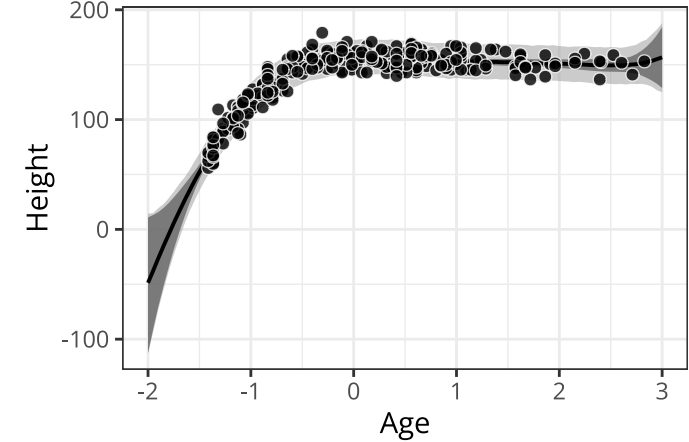
Predictions of mod3.4



Predictions of mod3.5



Predictions of mod3.6



Solution - Question 3

```
1 # prédictions moyennées sur les 4 modèles
2 averaged_predictions_mu <- pp_average(
3   mod3.1, mod3.2, mod3.3, mod3.4, mod3.5, mod3.6,
4   weights = "waic",
5   method = "fitted",
6   newdata = age_seq
7 ) %>%
8   as.data.frame() %>%
9   bind_cols(age_seq)
10
11 # prédictions moyennées sur les 4 modèles
12 averaged_predictions_age <- pp_average(
13   mod3.1, mod3.2, mod3.3, mod3.4, mod3.5, mod3.6,
14   weights = "waic",
15   method = "predict",
16   newdata = age_seq
17 ) %>%
18   as.data.frame() %>%
19   bind_cols(age_seq)
```

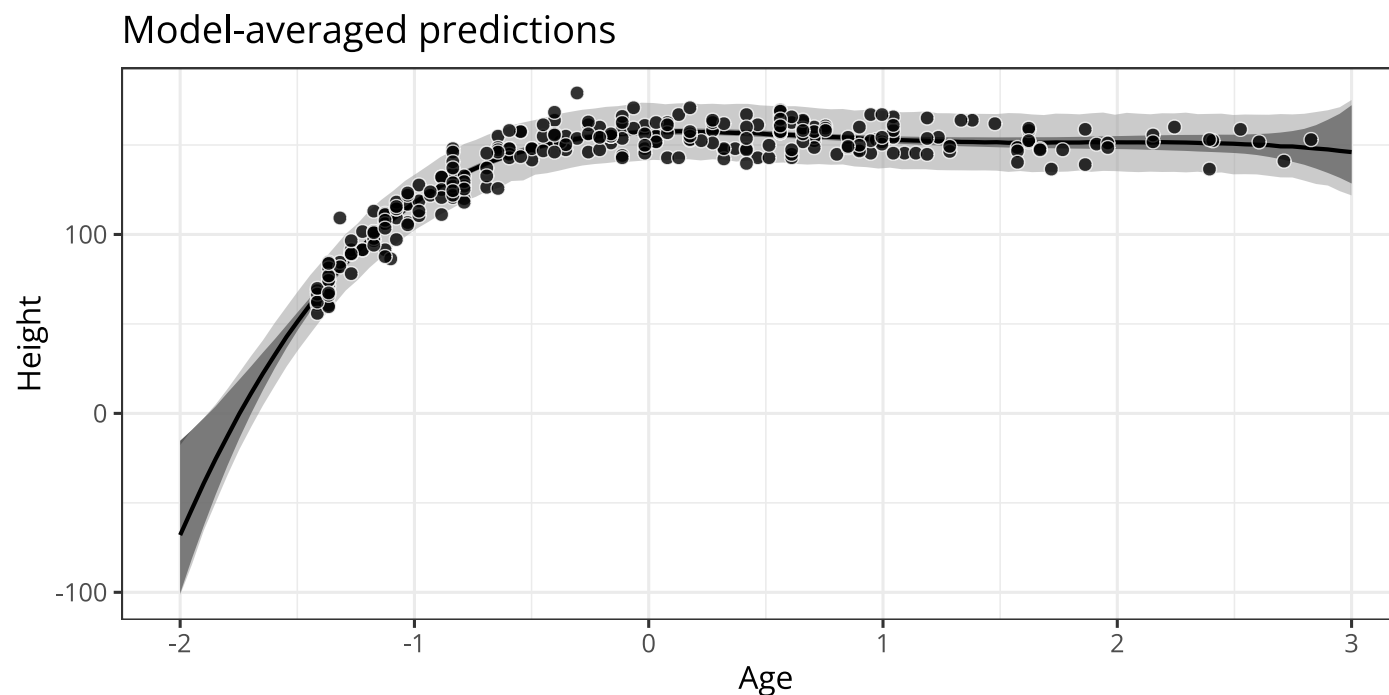


Solution - Question 3

```

1 d1 %>%
2   ggplot(aes(x = age, y = height) ) +
3   geom_ribbon(
4     data = averaged_predictions_mu, aes(x = age, ymin = Q2.5, ymax = Q97.5),
5     alpha = 0.8, inherit.aes = FALSE
6   ) +
7   geom_smooth(
8     data = averaged_predictions_age, aes(y = Estimate, ymin = Q2.5, ymax = Q97.5),
9     stat = "identity", color = "black", alpha = 0.5, size = 1
10  ) +
11  geom_point(colour = "white", fill = "black", pch = 21, size = 3, alpha = 0.8) +
12  labs(x = "Age", y = "Height", title = "Model-averaged predictions")

```



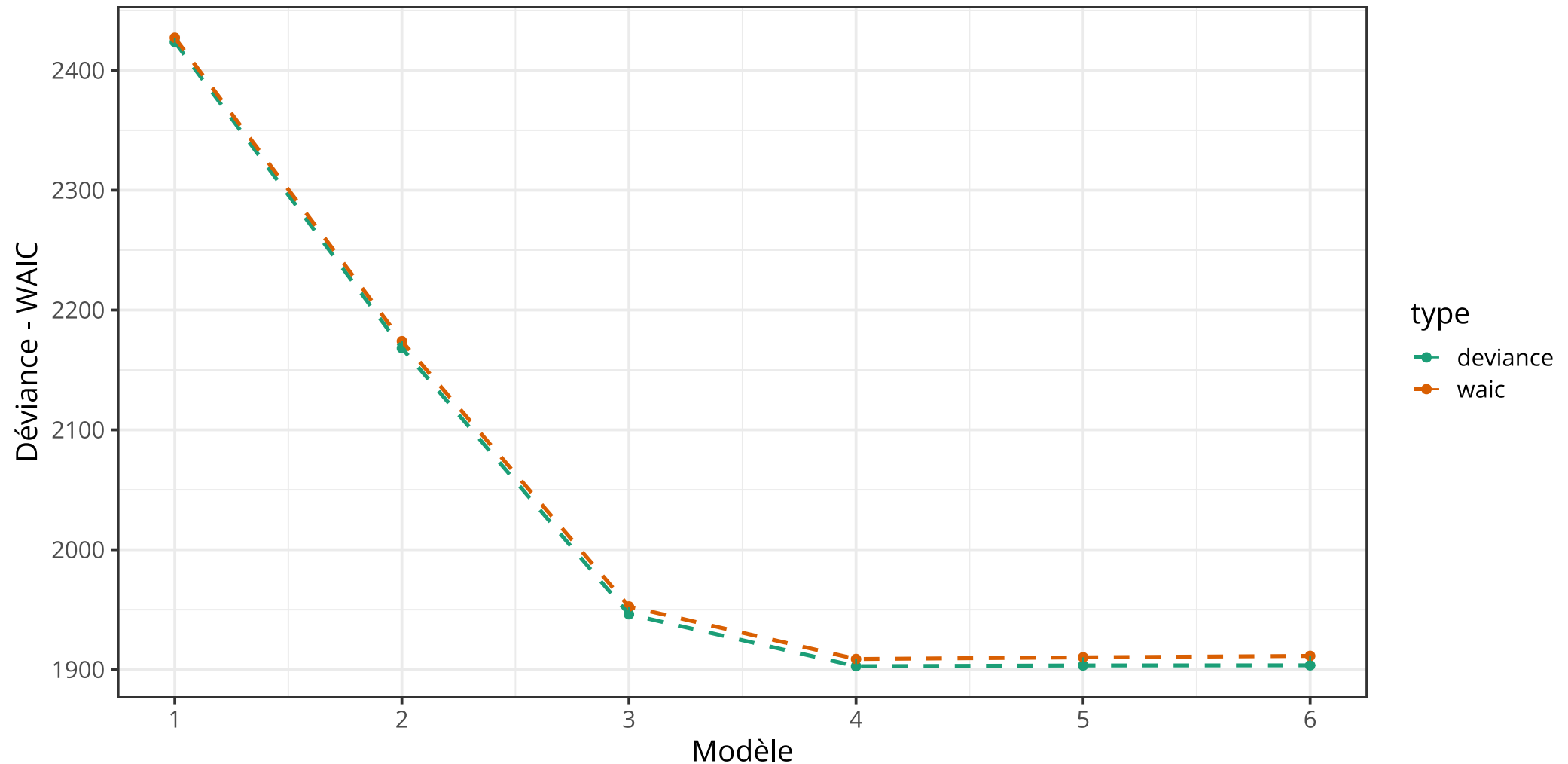
Solution - Question 4

```
1 # calcul de la log-vraisemblance (log-likelihood) du modèle mod3.1
2 # chaque ligne est une itération et chaque colonne une observation
3 log_lik_mod3.1 <- log_lik(mod3.1)
4
5 # NB : la déviance possède également une distribution dans le monde bayésien...
6 dev.mod3.1 <- mean(-2 * rowSums(log_lik_mod3.1) )
7
8 # calcul de la log-vraisemblance (log-likelihood) du modèle mod3.2
9 dev.mod3.2 <- mean(-2 * rowSums(log_lik(mod3.2) ) )
10
11 # calcul de la log-vraisemblance (log-likelihood) du modèle mod3.3
12 dev.mod3.3 <- mean(-2 * rowSums(log_lik(mod3.3) ) )
13
14 # calcul de la log-vraisemblance (log-likelihood) du modèle mod3.4
15 dev.mod3.4 <- mean(-2 * rowSums(log_lik(mod3.4) ) )
16
17 # calcul de la log-vraisemblance (log-likelihood) du modèle mod3.5
18 dev.mod3.5 <- mean(-2 * rowSums(log_lik(mod3.5) ) )
19
20 # calcul de la log-vraisemblance (log-likelihood) du modèle mod3.6
21 dev.mod3.6 <- mean(-2 * rowSums(log_lik(mod3.6) ) )
```



Solution - Question 4

```
1 deviances <- c(dev.mod3.1, dev.mod3.2, dev.mod3.3, dev.mod3.4, dev.mod3.5, dev.mod3.6)
2 comparison <- mod_comp %>% data.frame %>% select(waic) %>% rownames_to_column()
3 waics <- comparison %>% arrange(rowname) %>% pull(waic)
```



Références

- Burnham, K. P., & Anderson, D. R. (2004). Multimodel inference: Understanding AIC and BIC in model selection. *Sociological Methods & Research*, 33(2), 261–304. <https://doi.org/10.1177/0049124104268644>
- Burnham, K. P., & Anderson, D. R. (2002). *Model selection and multimodel inference: A practical information-theoretic approach* (2nd ed). Springer.

