

# 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



# Introduction

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

$$\mu_i = \alpha + \beta_1 \times x_{1i} + \beta_2 \times x_{2i}$$

Le modèle linéaire Gaussien qu'on a vu aux Cours n°03 et n°04 est caractérisé par un ensemble de postulats, entre autres choses :

- Les résidus sont distribués selon une loi Normale.
- La variance de cette distribution Normale est constante (postulat d'homogénéité de la variance).
- Les prédicteurs agissent sur la moyenne de cette distribution.
- La moyenne suit un modèle linéaire ou multi-linéaire.



# Introduction

Cette modélisation (le choix d'une distribution Normale) induit plusieurs contraintes, par exemple :

- Les données observées sont définies sur un espace continu.
- Cet espace n'est pas borné.

Comment modéliser certaines données qui ne respectent pas ces contraintes ? Par exemple, la proportion de bonnes réponses à un test (bornée entre 0 et 1), un temps de réponse (restreint à des valeurs positives et souvent distribué de manière approximativement log-normale), un nombre d'avalanches...



# Introduction

Nous avons déjà rencontré un modèle différent : le modèle Beta-Binomial (cf. Cours n°02).

$$y_i \sim \text{Binomial}(n, p)$$

$$p \sim \text{Beta}(a, b)$$

- Les données observées sont binaires (e.g., 0 vs. 1) ou le résultat d'une somme d'observations binaires (e.g., une proportion).
- La probabilité de succès (obtenir 1) a priori se caractérise par une distribution Beta.
- La probabilité de succès (obtenir 1) ne dépend d'aucun prédicteur.



# Introduction

Cette modélisation induit deux contraintes :

- Les données observées sont définies sur un espace discret.
- Cet espace est borné.

Comment pourrait-on ajouter des prédicteurs à ce modèle ?



# Modèle linéaire généralisé

$$y_i \sim \text{Binomial}(n, p_i)$$
$$f(p_i) = \alpha + \beta \times x_i$$

Objectifs :

- Rendre compte de données discrètes (e.g., échec/succès) générées par un processus unique.
- Introduire des prédicteurs dans le modèle.

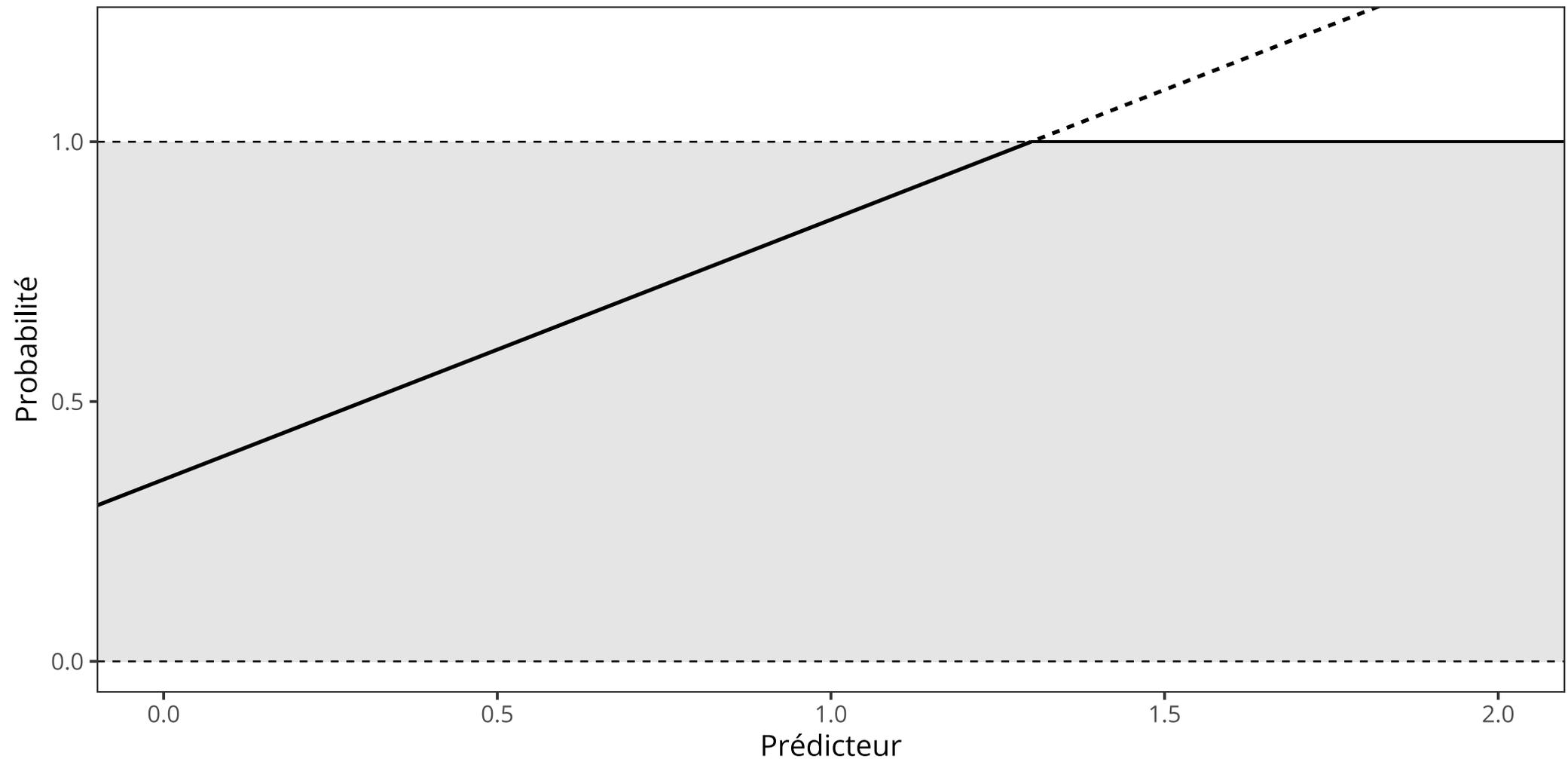
Deux changements par rapport au modèle Gaussien :

- L'utilisation d'une distribution de probabilité Binomiale.
- Le modèle linéaire ne sert plus à décrire directement un des paramètres de la distribution, mais une fonction de ce paramètre (on peut aussi considérer que le modèle Gaussien était formulé avec une fonction identité).



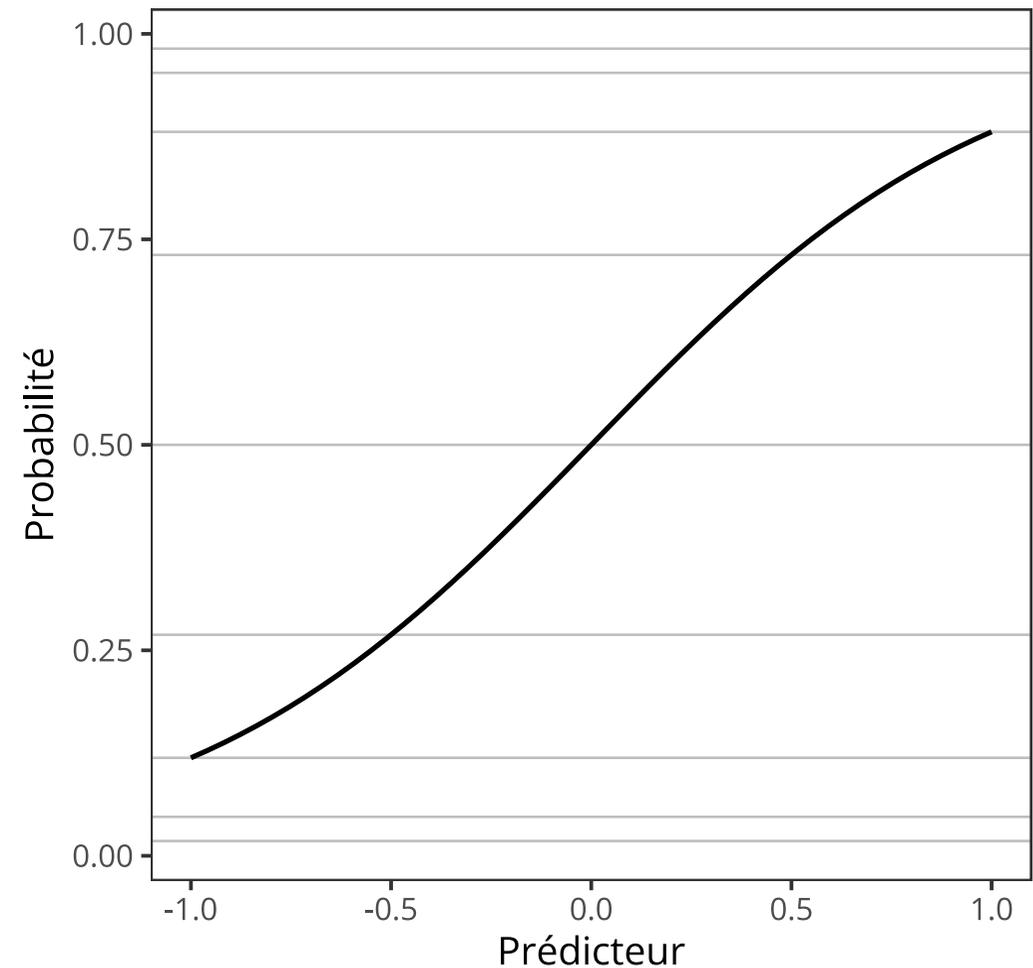
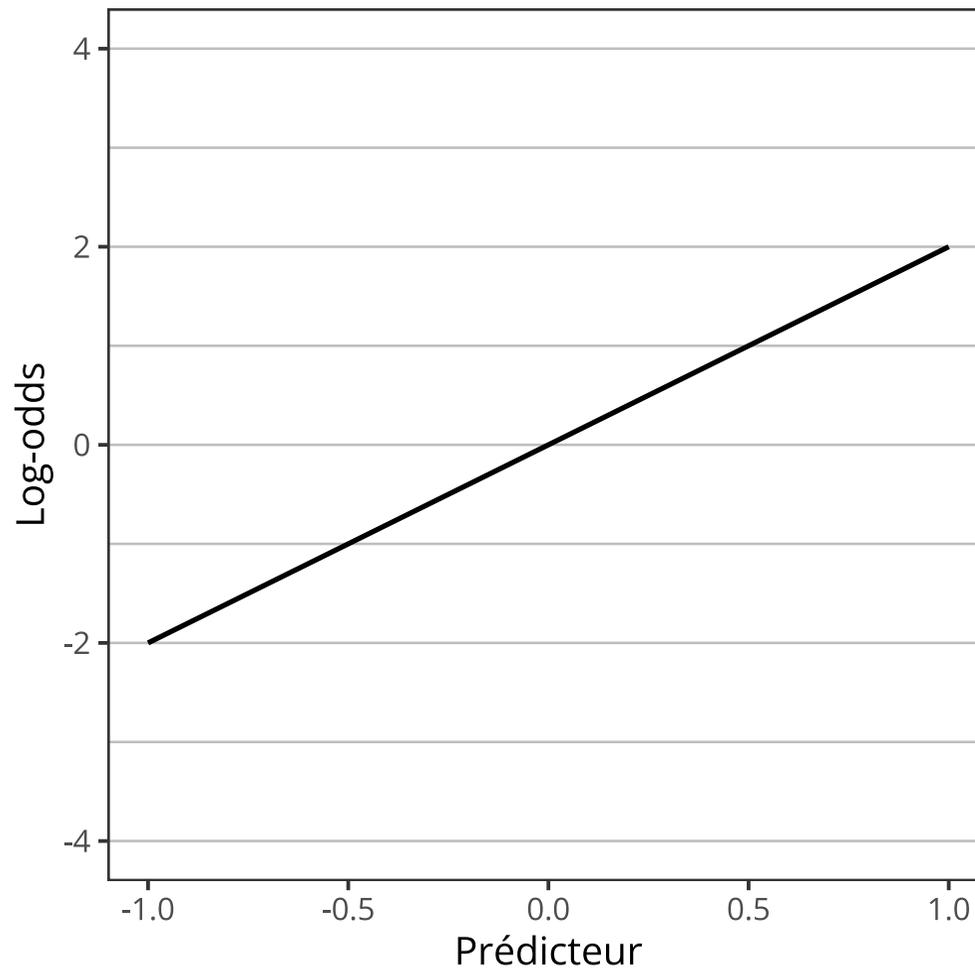
# Fonction de lien

Les fonctions de lien ont pour tâche de mettre en correspondance l'espace d'un modèle linéaire (non borné) avec l'espace d'un paramètre potentiellement borné comme une probabilité, définie sur l'intervalle  $[0, 1]$ .



# Fonction de lien

Les fonctions de lien ont pour tâche de mettre en correspondance l'espace d'un modèle linéaire (non borné) avec l'espace d'un paramètre potentiellement borné comme une probabilité, définie sur l'intervalle  $[0, 1]$ .



# Régression logistique

La fonction logit du GLM binomial (on parle de “log-odds”) :

$$\text{logit}(p_i) = \log\left(\frac{p_i}{1 - p_i}\right)$$

La cote d'un évènement (“odds” en anglais) est le ratio entre la probabilité que l'évènement se produise et la probabilité qu'il ne se produise pas. Le logarithme de cette cote est prédit par un modèle linéaire.

$$\log\left(\frac{p_i}{1 - p_i}\right) = \alpha + \beta \times x_i$$

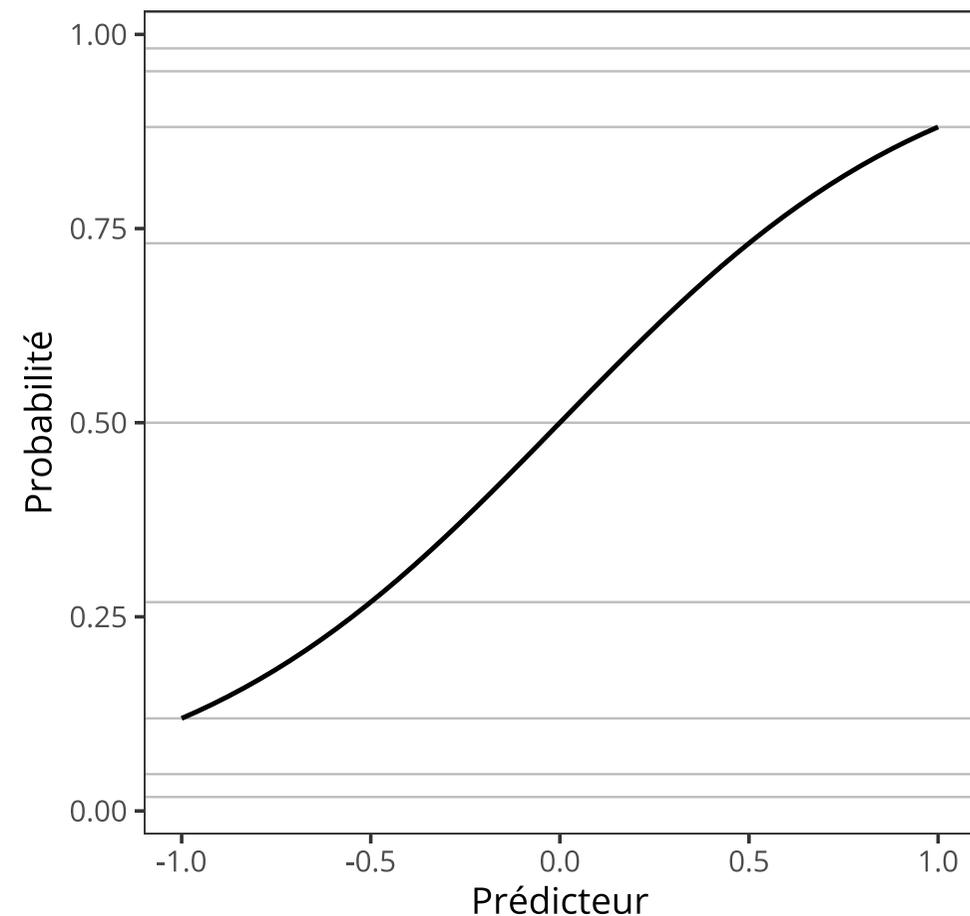
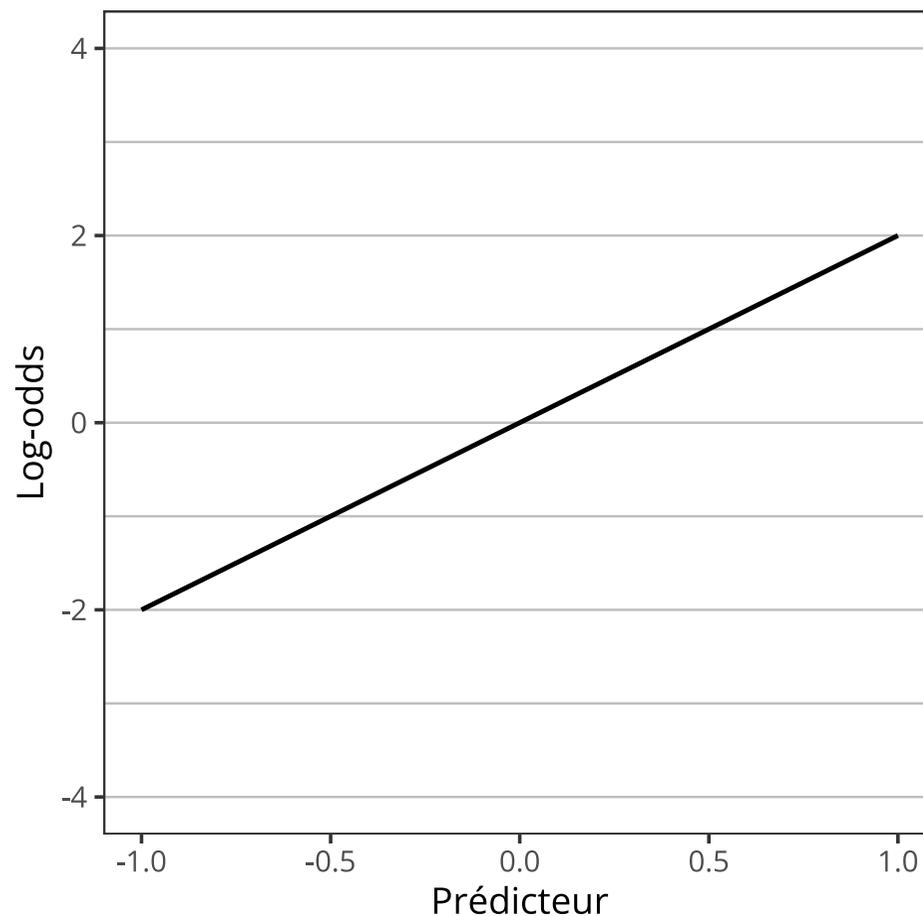
Pour retrouver la probabilité d'un évènement, il faut utiliser la fonction de **lien inverse**, la fonction **logistique** (ou logit-inverse) :

$$p_i = \frac{\exp(\alpha + \beta \times x_i)}{1 + \exp(\alpha + \beta \times x_i)}$$



# Complications induites par la fonction de lien

Ces fonctions de lien posent des problèmes d'interprétation : Un changement d'une unité sur un prédicteur n'a plus un effet constant sur la probabilité mais la modifie plus ou moins en fonction de son éloignement à l'origine. Quand  $x = 0$ , une augmentation d'une demi-unité (i.e.,  $\Delta x = 0.5$ ) se traduit par une augmentation de la probabilité de 0.25. Puis, chaque augmentation d'une demi-unité se traduit par une augmentation de  $p$  de plus en plus petite...



# Complications induites par la fonction de lien

Deuxième complication : cette fonction de lien “force” chaque prédicteur à interagir avec lui même et à interagir avec TOUS les autres prédicteurs, même si les interactions ne sont pas explicites...

Dans un modèle Gaussien, le taux de changement de  $y$  en fonction de  $x$  est donné par  $\partial(\alpha + \beta x) / \partial x = \beta$  et ne dépend pas de  $x$  (i.e.,  $\beta$  est constant).

Dans un GLM binomial (avec une fonction de lien logit), la probabilité d'un évènement est donnée par la fonction logistique :

$$p_i = \frac{\exp(\alpha + \beta \times x_i)}{1 + \exp(\alpha + \beta \times x_i)}$$

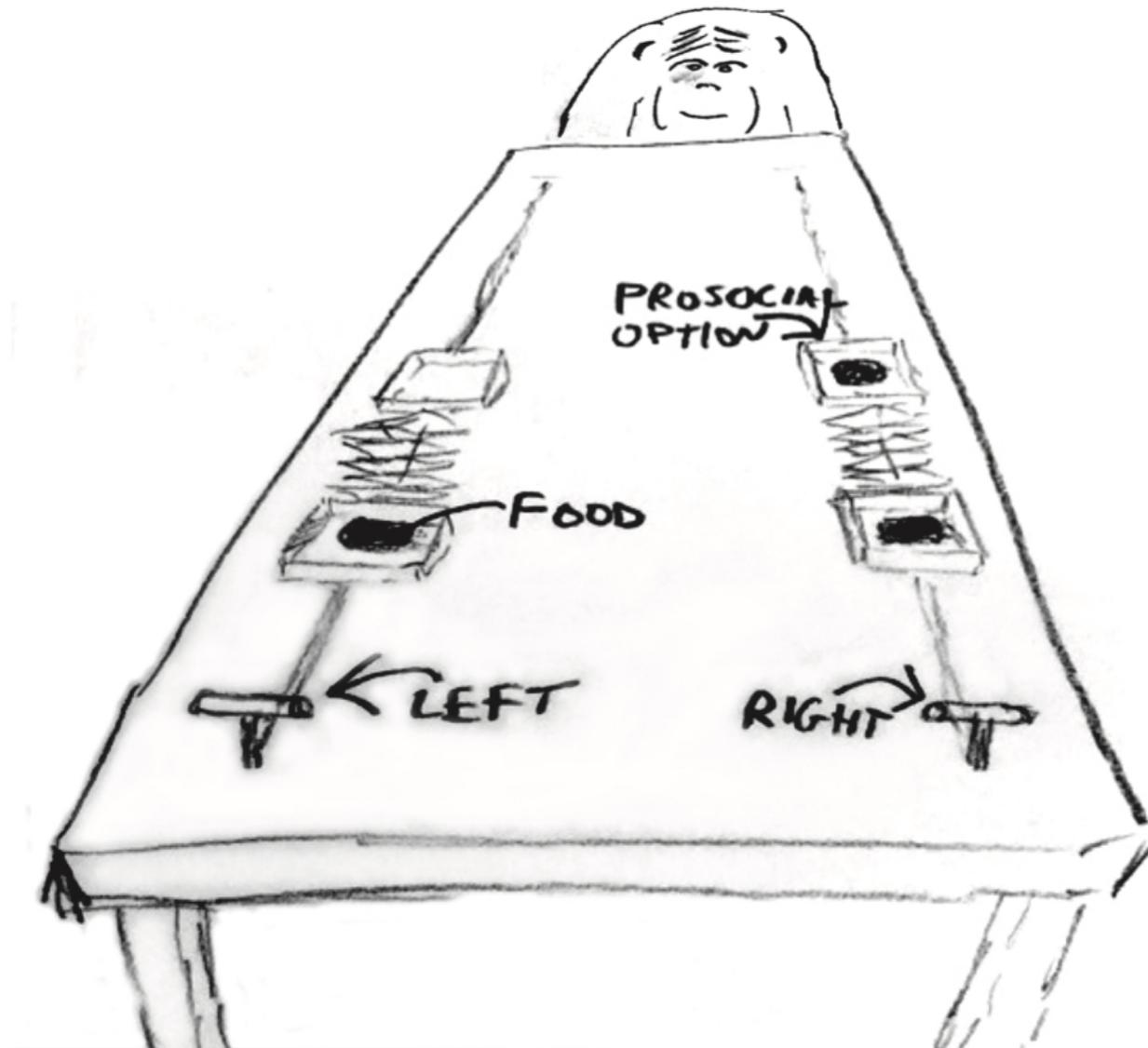
Et le taux de changement de  $p$  en fonction du prédicteur  $x$  est donné par :

$$\frac{\partial p}{\partial x} = \frac{\beta}{2(1 + \cosh(\alpha + \beta \times x))}$$

On voit que la variation sur  $p$  due au prédicteur  $x$  est fonction du prédicteur  $x$ , et dépend également de la valeur de  $\alpha$ ... !



# Exemple de régression logistique : La prosocialité chez le chimpanzé



# Régression logistique

```

1 library(tidyverse)
2 library(imsb)
3
4 df1 <- open_data(chimpanzees)
5 str(df1)

```

```

'data.frame':  504 obs. of  8 variables:
 $ actor      : int  1 1 1 1 1 1 1 1 1 1 ...
 $ recipient  : int  NA ...
 $ condition  : int  0 0 0 0 0 0 0 0 0 0 ...
 $ block      : int  1 1 1 1 1 1 2 2 2 2 ...
 $ trial      : int  2 4 6 8 10 12 14 16 18 20 ...
 $ prosoc_left : int  0 0 1 0 1 1 1 1 0 0 ...
 $ chose_prosoc: int  1 0 0 1 1 1 0 0 1 1 ...
 $ pulled_left : int  0 1 0 0 1 1 0 0 0 0 ...

```

- **pulled\_left** : 1 lorsque le chimpanzé pousse le levier gauche, 0 sinon.
- **prosoc\_left** : 1 lorsque le levier gauche est associé à l'option prosociale, 0 sinon.
- **condition** : 1 lorsqu'un partenaire est présent, 0 sinon.



# Régression logistique

## Le problème

On cherche à savoir si la présence d'un singe partenaire incite le chimpanzé à appuyer sur le levier prosocial, c'est à dire l'option qui donne de la nourriture aux deux individus. Autrement dit, est-ce qu'il existe une interaction entre l'effet de la latéralité et l'effet de la présence d'un autre chimpanzé sur la probabilité d'actionner le levier gauche.

## Les variables

- Observations (**pulled\_left**) : Ce sont des variables de Bernoulli. Elles prennent comme valeur 0/1.
- Prédicteur (**prosoc\_left**) : Est-ce que les deux plats sont sur la gauche ou sur la droite ?
- Prédicteur (**condition**) : Est-ce qu'un partenaire est présent ?



# Régression logistique

$$L_i \sim \text{Binomial}(1, p_i)$$

(équivalent à)  $L_i \sim \text{Bernoulli}(p_i)$

$$\text{logit}(p_i) = \alpha$$

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

Modèle mathématique sans prédicteur. Comment choisir une valeur pour  $\omega$ ... ?



# Prior predictive check

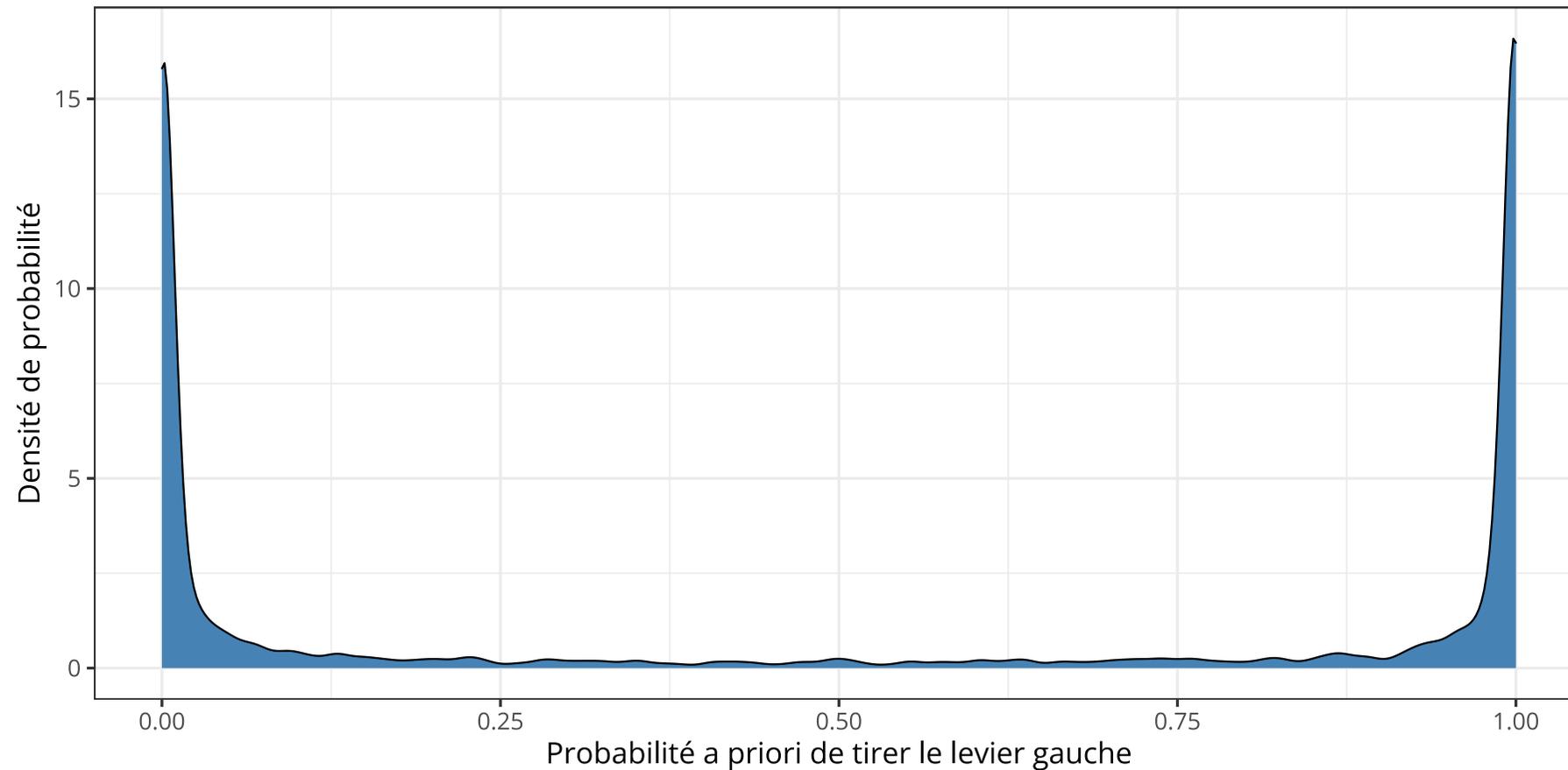
On écrit le modèle précédent avec `brms` et on échantillonne à partir du prior afin de vérifier que les prédictions du modèle (sur la base du prior seul) correspondent à nos attentes.

```
1 library(brms)
2
3 mod1.1 <- brm(
4   # "trials" permet de définir le nombre d'essais (i.e., n)
5   formula = pulled_left | trials(1) ~ 1,
6   family = binomial(),
7   prior = prior(normal(0, 10), class = Intercept),
8   data = df1,
9   # on veut récupérer les échantillons issus du prior
10  sample_prior = "yes"
11 )
```

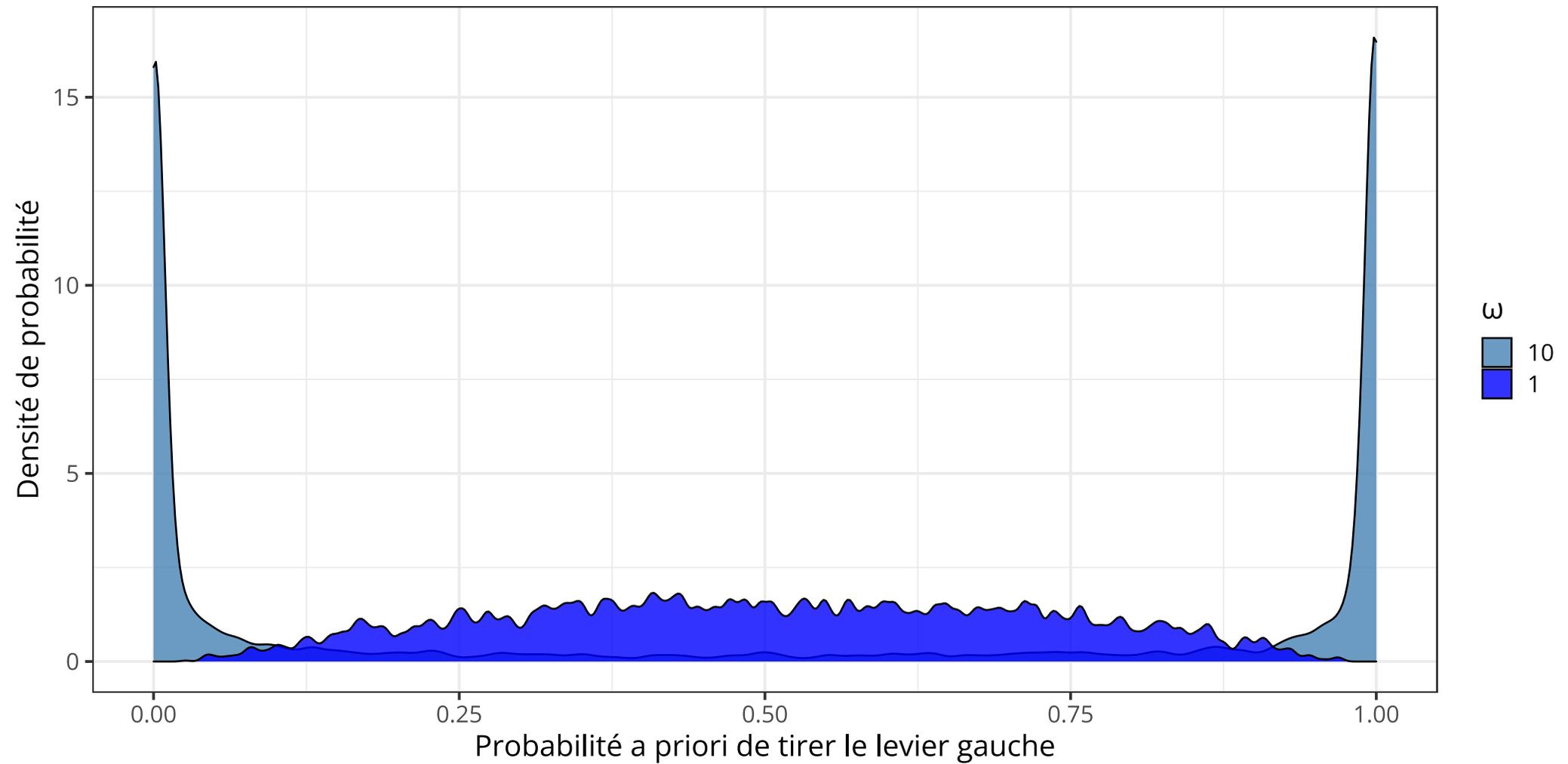


# Prior predictive check

```
1 # récupère les échantillons (sur la base) du prior
2 prior_draws(x = mod1.1) %>%
3   # applique la fonction de lien inverse
4   mutate(p = brms::inv_logit_scaled(Intercept) ) %>%
5   ggplot(aes(x = p) ) +
6   geom_density(fill = "steelblue", adjust = 0.1) +
7   labs(x = "Probabilité a priori de tirer le levier gauche", y = "Densité de probabilité")
```



# Prior predictive check



# Régression logistique

L'intercept s'interprète dans l'espace des log-odds... pour l'interpréter comme une probabilité, il faut appliquer la fonction de lien inverse. On peut utiliser la fonction `brms::inv_logit_scaled()` ou la fonction `plogis()`.

```
1 fixed_effects <- fixef(mod1.2) # effets fixes (i.e., que l'intercept ici)
2 plogis(fixed_effects) # fonction de lien inverse
```

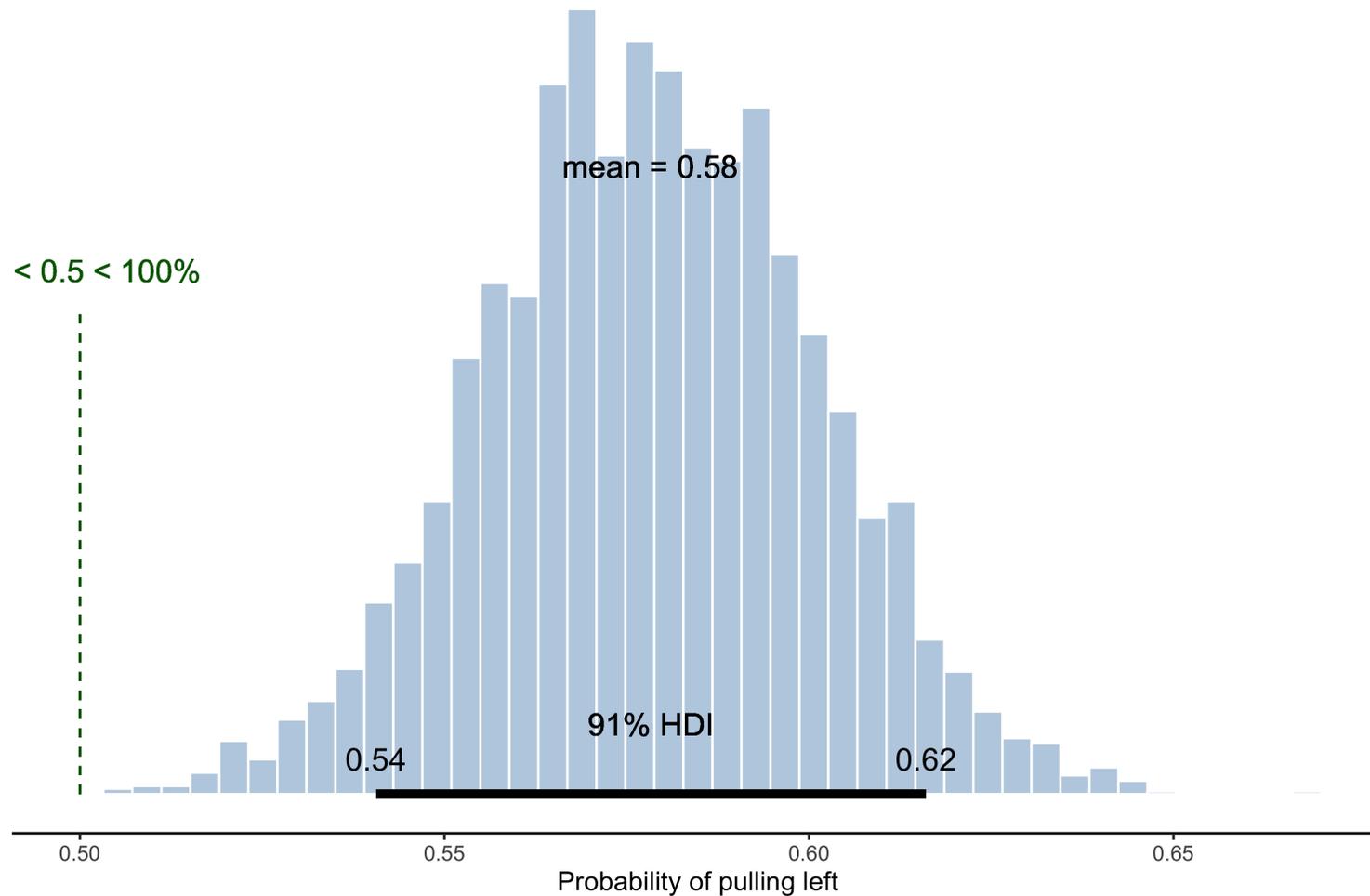
	Estimate	Est.Error	Q2.5	Q97.5
Intercept	0.5783937	0.5230697	0.5337498	0.6220639

En moyenne (sans considérer les prédicteurs), il semblerait que les chimpanzés aient légèrement plus tendance à appuyer sur le levier gauche que sur le levier droit...



# Régression logistique

```
1 post <- as_draws_df(x = mod1.2) # récupère les échantillons du posterior
2 intercept_samples <- plogis(post$b_Intercept) # échantillons pour l'intercept
3
4 posterior_plot(samples = intercept_samples, compval = 0.5) + labs(x = "Probability of pulling left")
```



# Régression logistique

Et si on ajoutait des prédicteurs... comment choisir une valeur pour  $\omega$  ?

$$L_i \sim \text{Binomial}(1, p_i)$$

$$\text{logit}(p_i) = \alpha + \beta_P P_i + \beta_C C_i + \beta_{PC} P_i C_i$$

$$\alpha \sim \text{Normal}(0, 1)$$

$$\beta_P, \beta_C, \beta_{PC} \sim \text{Normal}(0, \omega)$$

- $L_i$  indique si le singe a poussé le levier gauche (`pulled_left`).
- $P_i$  indique si le coté gauche correspond au coté prosocial.
- $C_i$  indique la présence d'un partenaire.



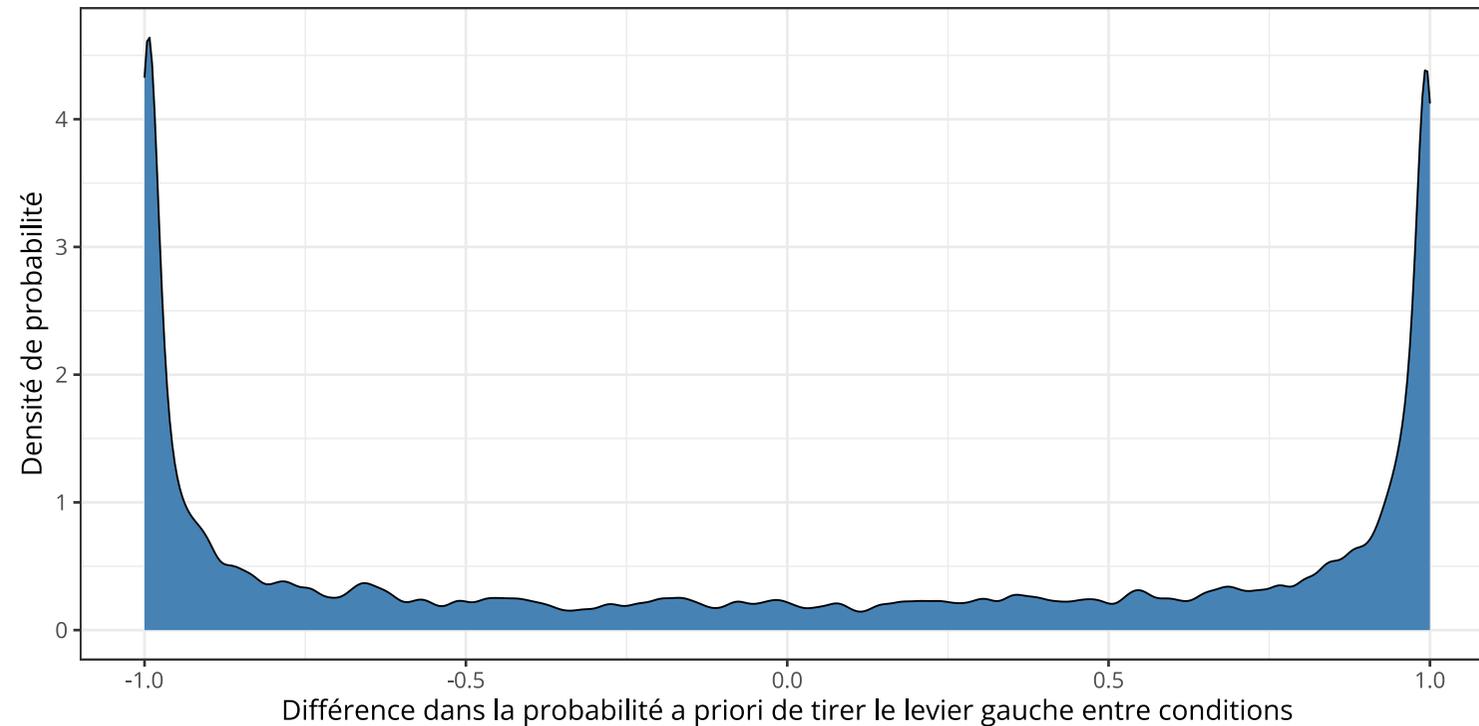
# Régression logistique

```
1 # recoding predictors
2 df1 <- df1 %>%
3   mutate(
4     prosoc_left = ifelse(prosoc_left == 1, 0.5, -0.5),
5     condition = ifelse(condition == 1, 0.5, -0.5)
6   )
7
8 priors <- c(
9   prior(normal(0, 1), class = Intercept),
10  prior(normal(0, 10), class = b)
11 )
12
13 mod2.1 <- brm(
14   formula = pulled_left | trials(1) ~ 1 + prosoc_left * condition,
15   family = binomial,
16   prior = priors,
17   data = df1,
18   sample_prior = "yes"
19 )
```



# Prior predictive check

```
1 prior_draws(x = mod2.1) %>% # échantillons du prior
2   mutate(
3     condition1 = plogis(Intercept - 0.5 * b), # p dans condition 1
4     condition2 = plogis(Intercept + 0.5 * b) # p dans condition 0
5   ) %>%
6   ggplot(aes(x = condition2 - condition1) ) + # on plot la différence
7   geom_density(fill = "steelblue", adjust = 0.1) +
8   labs(
9     x = "Différence dans la probabilité a priori de tirer le levier gauche entre conditions",
10    y = "Densité de probabilité"
11  )
```



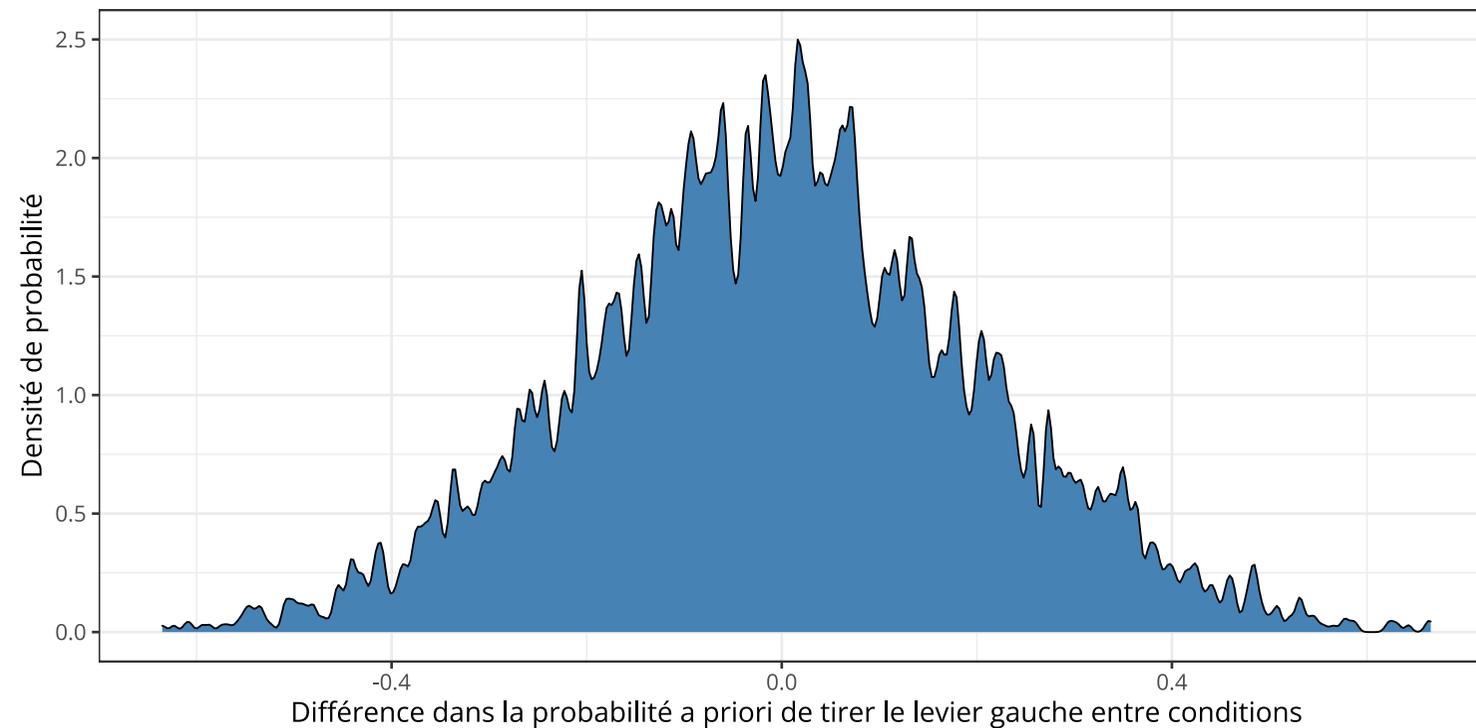
# Régression logistique

```
1 priors <- c(  
2   prior(normal(0, 1), class = Intercept),  
3   prior(normal(0, 1), class = b)  
4 )  
5  
6 mod2.2 <- brm(  
7   formula = pulled_left | trials(1) ~ 1 + prosoc_left * condition,  
8   family = binomial,  
9   prior = priors,  
10  data = df1,  
11  sample_prior = "yes"  
12 )
```



# Prior predictive check

```
1 prior_draws(mod2.2) %>% # échantillons du prior
2   mutate(
3     condition1 = plogis(Intercept - 0.5 * b), # p dans condition 1
4     condition2 = plogis(Intercept + 0.5 * b) # p dans condition 0
5   ) %>%
6   ggplot(aes(x = condition2 - condition1) ) +
7   geom_density(fill = "steelblue", adjust = 0.1) +
8   labs(
9     x = "Différence dans la probabilité a priori de tirer le levier gauche entre conditions",
10    y = "Densité de probabilité"
11  )
```



# Régression logistique

```
1 summary(mod2.2)
```

```
Family: binomial
Links: mu = logit
Formula: pulled_left | trials(1) ~ 1 + prosoc_left * condition
Data: df1 (Number of observations: 504)
Draws: 4 chains, each with iter = 2000; warmup = 1000; thin = 1;
       total post-warmup draws = 4000

Population-Level Effects:
              Estimate Est.Error 1-95% CI u-95% CI Rhat Bulk_ESS
Intercept          0.33      0.09   0.16   0.50 1.00    4221
prosoc_left         0.54      0.18   0.18   0.89 1.00    4682
condition          -0.19      0.18  -0.54   0.15 1.00    5133
prosoc_left:condition 0.16      0.34  -0.49   0.82 1.00    4691

              Tail_ESS
Intercept          3014
prosoc_left         3254
condition          3263
prosoc_left:condition 3253

Draws were sampled using sampling(NUTS). For each parameter, Bulk_ESS
and Tail_ESS are effective sample size measures, and Rhat is the potential
```



# Effet relatif vs. Effet absolu

Le modèle linéaire ne prédit pas directement la probabilité mais le log-odds de la probabilité :

$$\text{logit}(p_i) = \log\left(\frac{p_i}{1 - p_i}\right) = \alpha + \beta \times x_i$$

On peut distinguer et interpréter deux types d'effets.

**Effet relatif** : L'effet relatif porte sur le logarithme du rapport des probabilités. Il indique une proportion de changement induit par le prédicteur sur **les chances** de succès (ou plutôt, sur la cote). Cela ne nous dit rien de la probabilité de l'évènement, dans l'absolu.

**Effet absolu** : Effet qui porte directement sur la probabilité d'un évènement. Il dépend de tous les paramètres du modèle et nous donne l'impact effectif d'un changement d'une unité d'un prédicteur (dans l'espace des probabilités).



# Effet relatif

Il s'agit d'une **proportion** de changement induit par le prédicteur sur le rapport des chances ou "cote" (odds). Illustration avec un modèle sans interaction.

$$\log\left(\frac{p_i}{1-p_i}\right) = \alpha + \beta x_i$$

$$\frac{p_i}{1-p_i} = \exp(\alpha + \beta x_i)$$

La cote proportionnelle  $q$  d'un évènement est le nombre par lequel la cote est multipliée lorsque  $x_i$  augmente d'une unité.

$$q = \frac{\exp(\alpha + \beta(x_i + 1))}{\exp(\alpha + \beta x_i)} = \frac{\exp(\alpha) \exp(\beta x_i) \exp(\beta)}{\exp(\alpha) \exp(\beta x_i)} = \exp(\beta)$$

Lorsque  $q = 2$  (par exemple), une augmentation de  $x_i$  d'une unité génère un doublement de la cote.



# Interprétation de l'effet relatif

L'effet relatif d'un paramètre **dépend également des autres paramètres**. Dans le modèle précédent, le prédicteur `prosoc_left` augmente le logarithme de la cote d'environ 0.54, ce qui se traduit par une augmentation de la cote de  $\exp(0.54) \approx 1.72$  soit une augmentation d'environ 72% de la cote.

Supposons que l'intercept  $\alpha = 4$ .

- La probabilité de pousser le levier sans autre considération est de  $\text{logit}^{-1}(4) \approx 0.98$ .
- En considérant l'effet de `prosoc_left`, on obtient  $\text{logit}^{-1}(4 + 0.54) \approx 0.99$ .

Une augmentation de 72% sur le log-odds se traduit par une augmentation de seulement 1% sur la probabilité effective... Les effets relatifs peuvent conduire à de mauvaises interprétations lorsqu'on ne considère pas l'échelle de la variable mesurée.

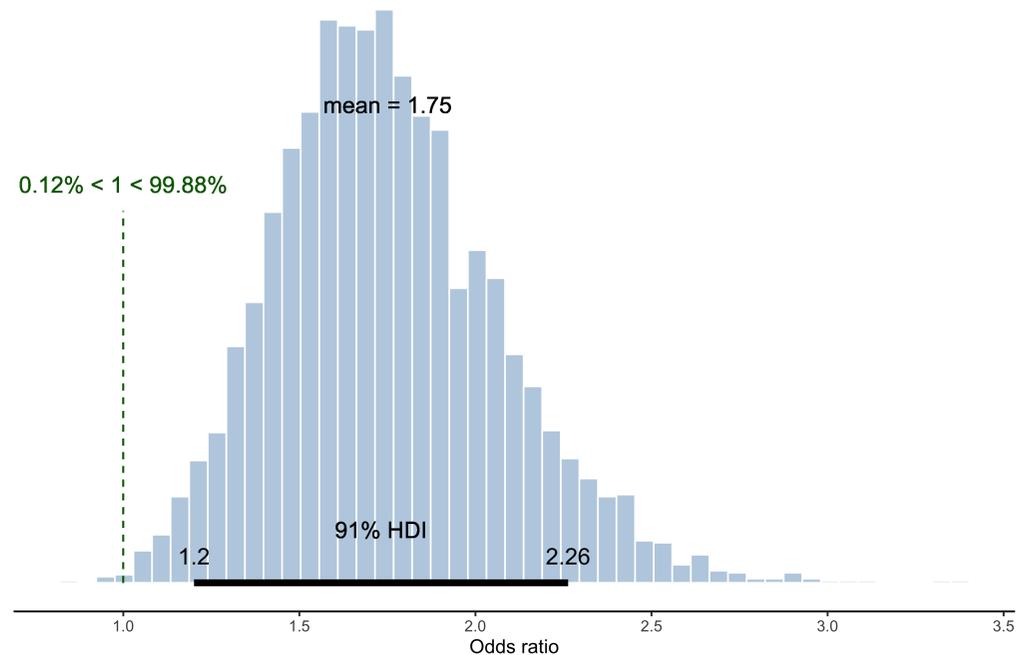


# Interprétation de l'effet relatif

```
1 fixef(mod2.2) # récupère les estimations des effets dits "fixes"
```

	Estimate	Est.Error	Q2.5	Q97.5
Intercept	0.3272441	0.09031635	0.1576418	0.5022809
prosoc_left	0.5435832	0.18088414	0.1819663	0.8936373
condition	-0.1905311	0.17859193	-0.5365269	0.1501853
prosoc_left:condition	0.1631846	0.33861972	-0.4929849	0.8198021

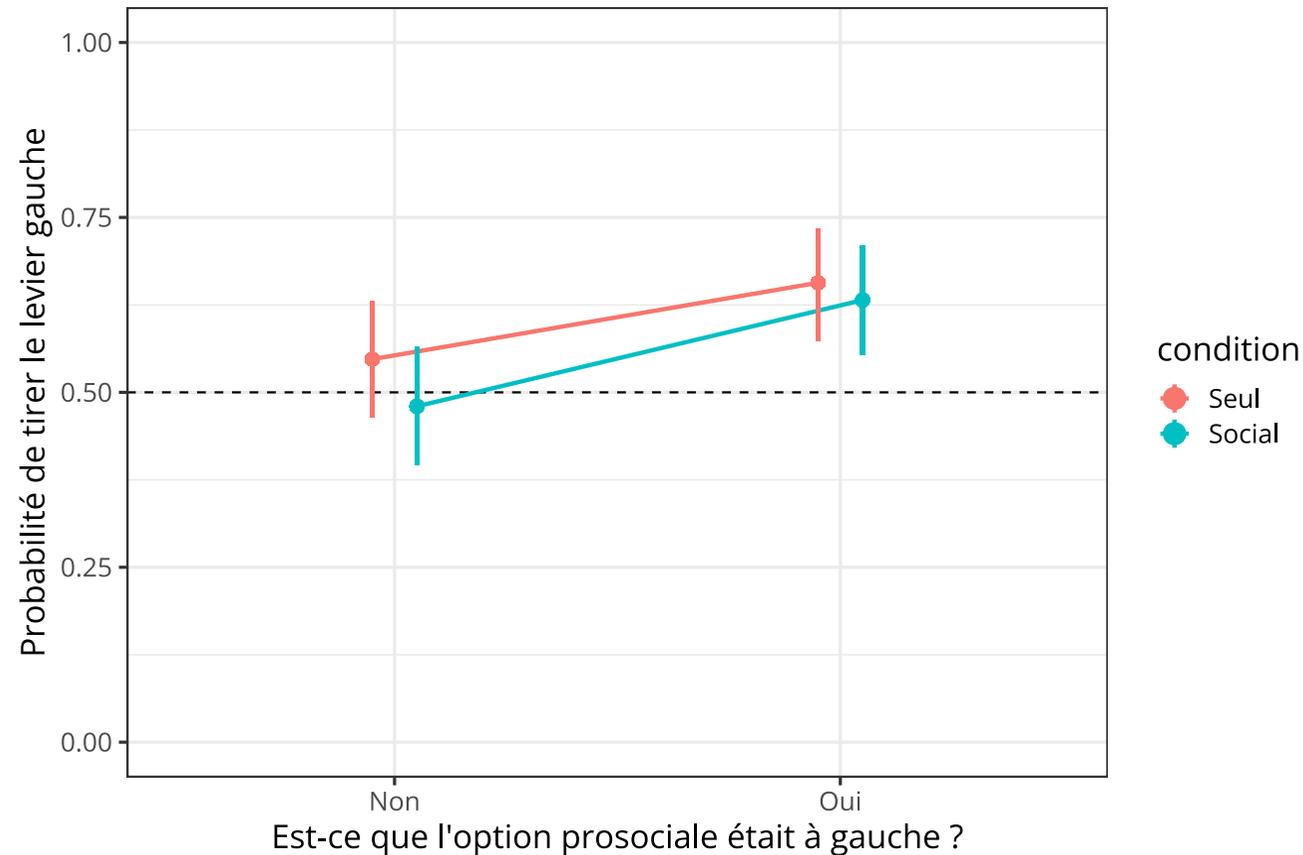
```
1 post <- as_draws_df(x = mod2.2) # échantillons du posterior
2 posterior_plot(samples = exp(post$b_prosoc_left), compval = 1) + labs(x = "Odds ratio")
```



# Effet absolu

L'effet absolu dépend de tous les paramètres du modèle et nous donne l'impact effectif d'un changement d'une unité d'un prédicteur (dans l'espace des probabilités).

```
1 model_predictions <- fitted(mod2.2) %>% # prédiction pour p (i.e., la probabilité)
2   data.frame() %>%
3   bind_cols(df1) %>%
4   mutate(condition = factor(condition), prosoc_left = factor(prosoc_left) )
```



# Régression binomiale agrégée

Ces données représentent le nombre de candidatures à l'université de Berkeley par sexe et par département. Chaque candidature est acceptée ou rejetée et les résultats sont agrégés par département et par sexe.

```
1 (df2 <- open_data(admission) )
```

```
# A tibble: 12 × 5
  dept gender admit reject applications
  <fct> <fct> <dbl> <dbl> <dbl>
1 A Male 512 313 825
2 A Female 89 19 108
3 B Male 353 207 560
4 B Female 17 8 25
5 C Male 120 205 325
6 C Female 202 391 593
7 D Male 138 279 417
8 D Female 131 244 375
9 E Male 53 138 191
10 E Female 94 299 393
11 F Male 22 351 373
12 F Female 24 317 341
```

Existe-t-il un biais de recrutement lié au sexe ?



# Régression binomiale agrégée

On va construire un modèle de la décision d'admission en prenant comme prédicteur le sexe du candidat.

$$\begin{aligned} \text{admit}_i &\sim \text{Binomial}(n_i, p_i) \\ \text{logit}(p_i) &= \alpha + \beta_m \times m_i \\ \alpha &\sim \text{Normal}(0, 1) \\ \beta_m &\sim \text{Normal}(0, 1) \end{aligned}$$

Les variables :

- $\text{admit}_i$  : Le nombre de candidatures acceptées (**admit**).
- $n_i$  : Le nombre total de candidatures (**applications**).
- $m_i$  : Le sexe du candidat (**1 = Male**).



# Régression binomiale agrégée

```
1 priors <- c(prior(normal(0, 1), class = Intercept) )
2
3 mod3 <- brm(
4   formula = admit | trials(applications) ~ 1,
5   family = binomial(link = "logit"),
6   prior = priors,
7   data = df2,
8   sample_prior = "yes"
9 )
```



# Régression binomiale agrégée

```
1 priors <- c(  
2   prior(normal(0, 1), class = Intercept),  
3   prior(normal(0, 1), class = b)  
4 )  
5  
6 # dummy-coding  
7 df2$male <- ifelse(df2$gender == "Male", 1, 0)  
8  
9 mod4 <- brm(  
10  formula = admit | trials(applications) ~ 1 + male,  
11  family = binomial(link = "logit"),  
12  prior = priors,  
13  data = df2,  
14  sample_prior = "yes"  
15 )
```



# Régression binomiale agrégée

```
1 summary(mod4)
```

```
Family: binomial
Links: mu = logit
Formula: admit | trials(applications) ~ 1 + male
Data: df2 (Number of observations: 12)
Draws: 4 chains, each with iter = 2000; warmup = 1000; thin = 1;
       total post-warmup draws = 4000

Population-Level Effects:
      Estimate Est.Error 1-95% CI u-95% CI Rhat Bulk_ESS Tail_ESS
Intercept  -0.83      0.05  -0.93  -0.74 1.00    2149    2196
male        0.61      0.06   0.49   0.73 1.00    2382    2528

Draws were sampled using sampling(NUTS). For each parameter, Bulk_ESS
and Tail_ESS are effective sample size measures, and Rhat is the potential
scale reduction factor on split chains (at convergence, Rhat = 1).
```

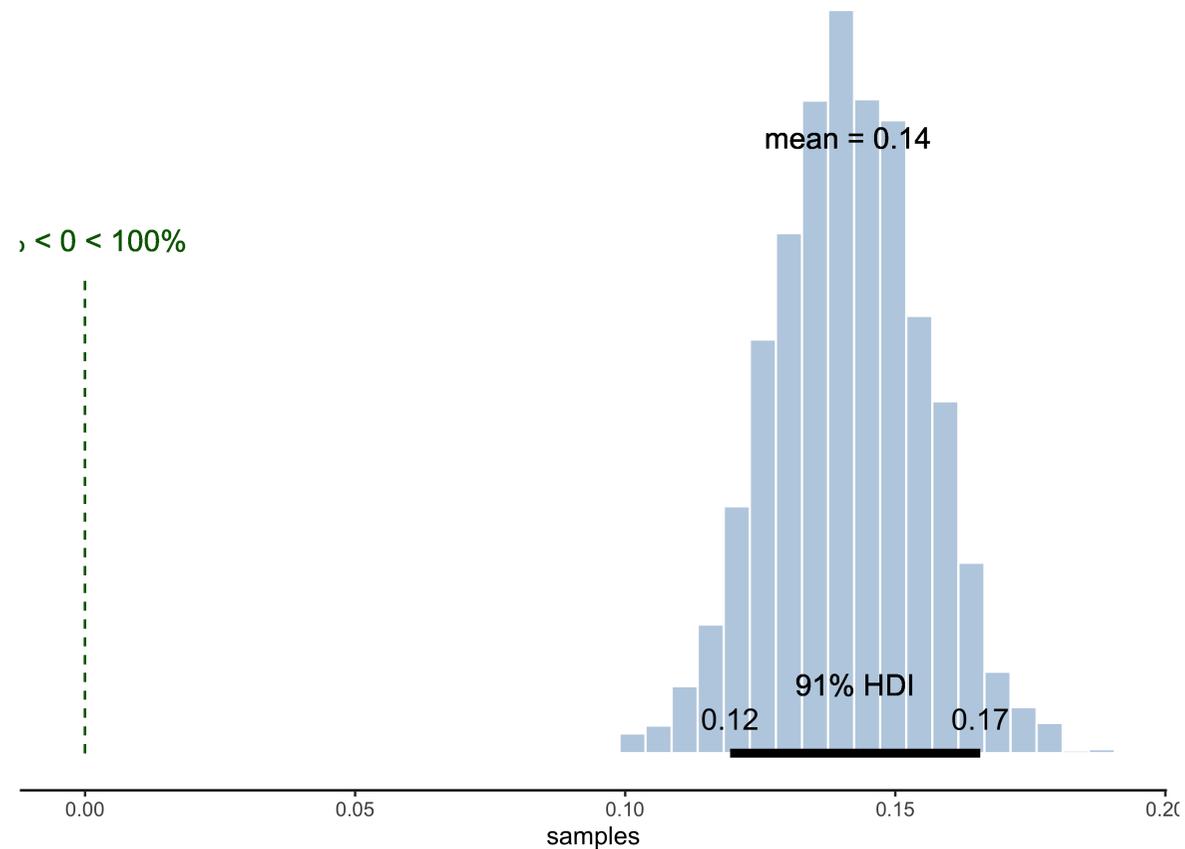
Être un homme semble être un avantage... ! Le rapport des cotes est de  $\exp(0.61) \approx 1.84$ .



# Régression binomiale agrégée

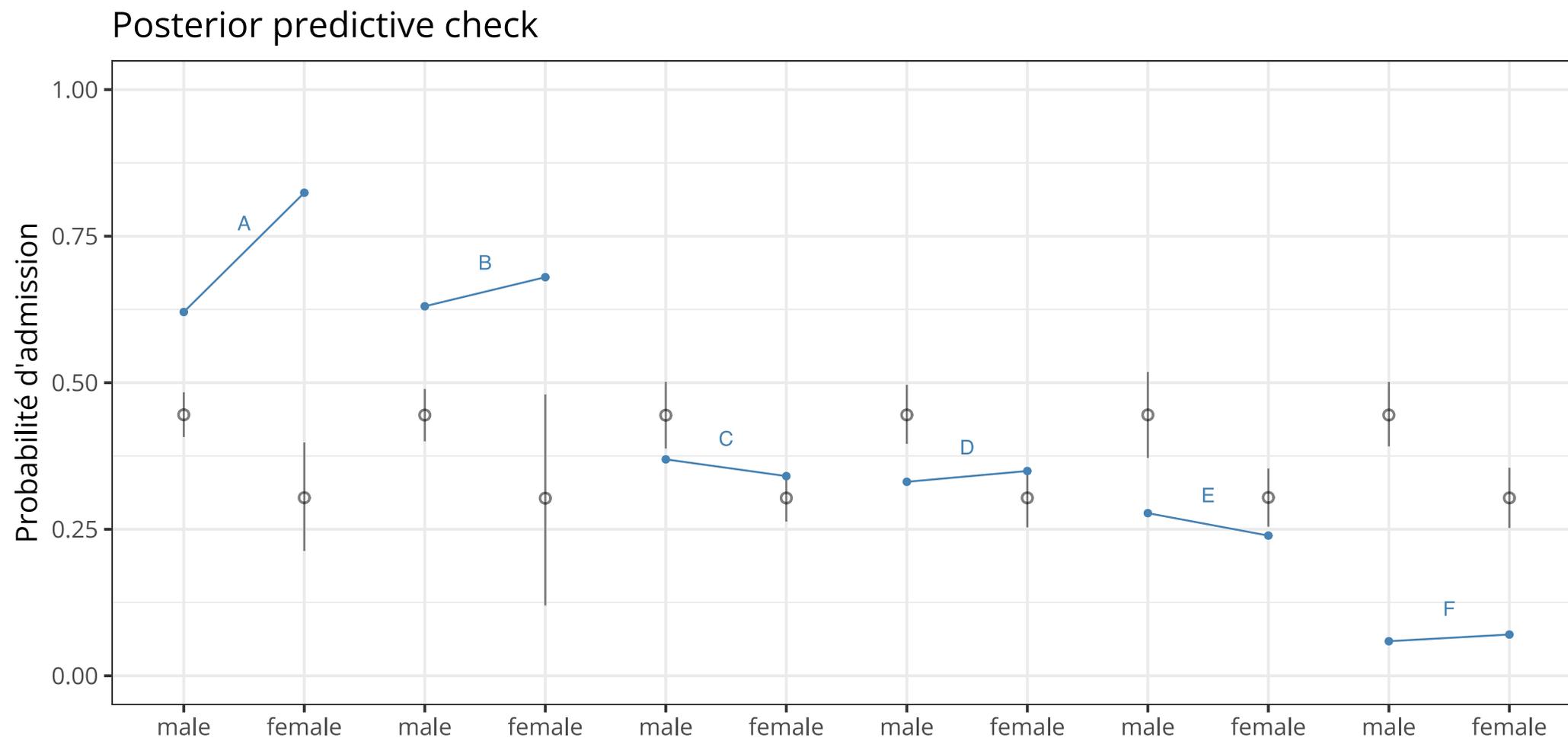
Calculons la différence de probabilité d'admission entre hommes et femmes.

```
1 post <- as_draws_df(x = mod4)
2 p.admit.male <- plogis(post$b_Intercept + post$b_male)
3 p.admit.female <- plogis(post$b_Intercept)
4 diff.admit <- p.admit.male - p.admit.female
5 posterior_plot(samples = diff.admit, compval = 0)
```



# Visualiser les prédictions du modèle

On examine les prédictions du modèle par département.



# Régression binomiale agrégée

Les prédictions du modèle sont très mauvaises... Il n'y a que deux départements pour lesquels les femmes ont de moins bonnes prédictions que les hommes (C et E) alors que le modèle prédit une probabilité d'admission plus basse pour tous les départements...

Le problème est double :

- Les hommes et les femmes ne postulent pas aux mêmes départements.
- Les départements n'ont pas tous les mêmes effectifs.

C'est le "paradoxe" de Simpson... remarques :

- La distribution postérieure seule n'aurait pas permis de détecter ce problème.
- C'est l'étude des prédictions du modèle qui nous a permis de mettre le doigt sur le problème...



# Régression binomiale agrégée

On construit donc un modèle de la décision d'admission en fonction du genre, au sein de chaque département.

$$\text{admit}_i \sim \text{Binomial}(n_i, p_i)$$

$$\text{logit}(p_i) = \alpha_{\text{dept}[i]} + \beta_m \times m_i$$

$$\alpha_{\text{dept}[i]} \sim \text{Normal}(0, 1)$$

$$\beta_m \sim \text{Normal}(0, 1)$$



# Régression binomiale agrégée

```
1 # modèle sans prédicteur
2 mod5 <- brm(
3   admit | trials(applications) ~ 0 + dept,
4   family = binomial(link = "logit"),
5   prior = prior(normal(0, 1), class = b),
6   data = df2
7 )
8
9 # modèle avec prédicteur
10 mod6 <- brm(
11   admit | trials(applications) ~ 0 + dept + male,
12   family = binomial(link = "logit"),
13   prior = prior(normal(0, 1), class = b),
14   data = df2
15 )
```



# Régression binomiale agrégée

```
1 summary(mod6)
```

```
Family: binomial
Links: mu = logit
Formula: admit | trials(applications) ~ 0 + dept + male
Data: df2 (Number of observations: 12)
Draws: 4 chains, each with iter = 2000; warmup = 1000; thin = 1;
total post-warmup draws = 4000
```

Population-Level Effects:

	Estimate	Est.Error	l-95% CI	u-95% CI	Rhat	Bulk_ESS	Tail_ESS
deptA	0.68	0.10	0.49	0.87	1.00	2128	2827
deptB	0.64	0.11	0.42	0.85	1.00	2449	2941
deptC	-0.58	0.07	-0.72	-0.43	1.00	3669	2813
deptD	-0.61	0.08	-0.77	-0.44	1.00	2987	2768
deptE	-1.05	0.10	-1.24	-0.85	1.00	4449	3062
deptF	-2.58	0.15	-2.88	-2.28	1.00	3729	2648
male	-0.10	0.08	-0.26	0.05	1.00	1763	2449

Draws were sampled using sampling(NUTS). For each parameter, Bulk\_ESS and Tail\_ESS are effective sample size measures, and Rhat is the potential scale reduction factor on split chains (at convergence, Rhat = 1).



# Régression binomiale agrégée

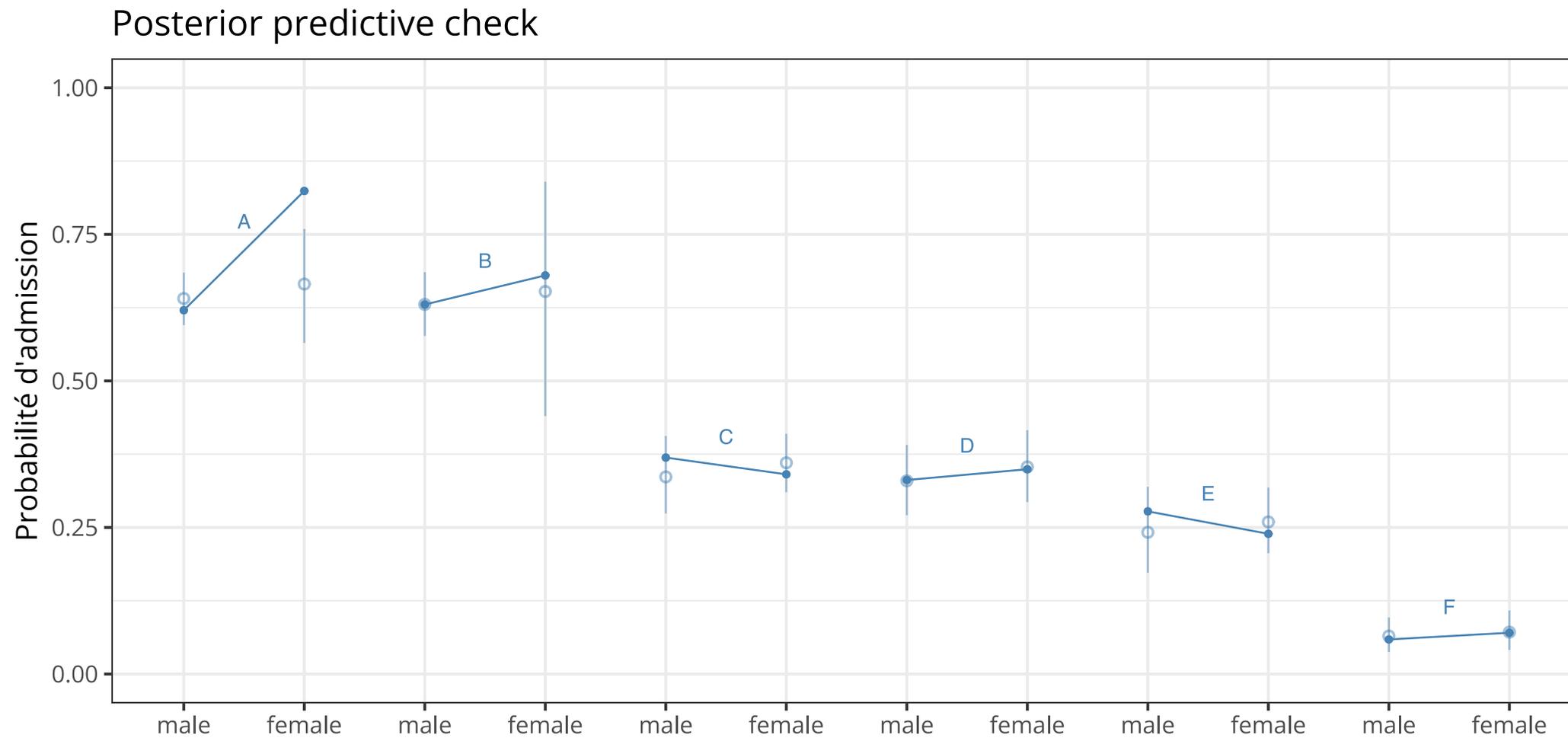
```
1 fixef(mod6)
```

	Estimate	Est.Error	Q2.5	Q97.5
deptA	0.6831831	0.09626582	0.4883134	0.86983987
deptB	0.6397385	0.11082315	0.4192897	0.85287772
deptC	-0.5774577	0.07333254	-0.7225987	-0.43446158
deptD	-0.6077253	0.08368156	-0.7683007	-0.44210784
deptE	-1.0479949	0.09914584	-1.2402143	-0.85198412
deptF	-2.5757293	0.15167036	-2.8810583	-2.28145169
male	-0.1038972	0.07921840	-0.2620235	0.05136751

Maintenant, la prédiction pour  $\beta_m$  va dans l'autre sens... Le rapport des cotes (odds ratio) est de  $\exp(-0.1) = 0.9$ , la cote (odds) des hommes est estimée à 90% de la cote des femmes.



# Régression binomiale agrégée



# Conclusions

Les hommes et les femmes ne postulent pas aux mêmes départements et les départements varient par leur probabilité d'admission. En l'occurrence, les femmes ont plus postulé aux départements E et F (avec une probabilité d'admission plus faible) et ont moins postulé aux départements A ou B, avec une probabilité d'admission plus haute.

Pour évaluer l'effet du sexe sur la probabilité d'admission, il faut donc se poser la question suivante : “Quelle est la différence de probabilité d'admission entre hommes et femmes **au sein de chaque département** ?” (plutôt que de manière générale).

Retenir que le modèle de régression peut être généralisé à différents modèles de génération des données (i.e., différentes distributions de probabilité, comme la distribution Normale, Binomiale, Poisson, etc) et que l'espace des paramètres peut être “connecté” à l'espace des prédicteurs (variables mesurées) grâce à des fonctions de lien (e.g., la fonction logarithme, exponentielle, logit, etc).

Retenir la distinction entre **effet relatif** (e.g., un changement de cote) et **effet absolu** (e.g., une différence de probabilité).



# Travaux pratiques - Absentéisme expérimental

Travailler avec des sujets humains implique un minimum de coopération réciproque. Mais ce n'est pas toujours le cas. Une partie non-négligeable des étudiants qui s'inscrivent pour passer des expériences de Psychologie ne se présentent pas le jour prévu... On a voulu estimer la **probabilité de présence d'un étudiant inscrit** en fonction de l'envoi (ou non) d'un mail de rappel (cet exemple est présenté en détails dans deux articles de blog, accessibles [ici](#), et [ici](#)).

```
1 df3 <- open_data(absence)
2 df3 %>% sample_frac %>% head(10)
```

	day	inscription	reminder	absence	presence	total
1	Tuesday	panel	yes	0	9	9
2	Thursday	doodle	no	3	11	14
3	Wednesday	doodle	no	6	11	17
4	Tuesday	doodle	no	4	10	14
5	Tuesday	doodle	yes	1	7	8
6	Friday	doodle	no	7	11	18
7	Friday	doodle	yes	0	2	2
8	Wednesday	doodle	yes	0	4	4
9	Monday	doodle	no	5	4	9
10	Monday	panel	yes	6	12	18



# Travaux pratiques

- **Quelle est la probabilité qu'un participant, qui s'est inscrit de son propre chef, vienne effectivement passer l'expérience ?**
- Quel est l'effet du rappel ?
- Quel est l'effet du mode d'inscription ?
- Quel est l'effet conjoint de ces deux prédicteurs ?



# Travaux pratiques

Écrire le modèle qui prédit la présence d'un participant sans prédicteur.

$$y_i \sim \text{Binomial}(n_i, p_i)$$
$$\text{logit}(p_i) = \alpha$$
$$\alpha \sim \text{Normal}(0, 1)$$



# Travaux pratiques

```

1 mod7 <- brm(
2   presence | trials(total) ~ 1,
3   family = binomial(link = "logit"),
4   prior = prior(normal(0, 1), class = Intercept),
5   data = df3,
6   # utilise tous les coeurs disponibles de la machine...
7   cores = parallel::detectCores()
8 )

```

```
1 fixef(mod7) # effet relatif (log de la cote)
```

	Estimate	Est.Error	Q2.5	Q97.5
Intercept	1.152125	0.1957073	0.7787238	1.553165

```
1 fixef(mod7) %>% plogis # effet absolu (probabilité de présence)
```

	Estimate	Est.Error	Q2.5	Q97.5
Intercept	0.7598988	0.5487713	0.685405	0.8253704



# Travaux pratiques

- Quelle est la probabilité qu'un participant, qui s'est inscrit de son propre chef, vienne effectivement passer l'expérience ?
- **Quel est l'effet du rappel ?**
- Quel est l'effet du mode d'inscription ?
- Quel est l'effet conjoint de ces deux prédicteurs ?



# Travaux pratiques

On commence par recoder en dummy variables `reminder` et `inscription`.

```

1 df3 <-
2   df3 %>%
3   mutate(
4     reminder = ifelse(reminder == "no", 0, 1),
5     inscription = ifelse(inscription == "panel", 0, 1)
6   )
7
8 head(df3, n = 10)

```

	day	inscription	reminder	absence	presence	total
1	Friday	1	0	7	11	18
2	Friday	1	1	0	2	2
3	Friday	0	1	0	10	10
4	Monday	1	0	5	4	9
5	Monday	1	1	2	6	8
6	Monday	0	1	6	12	18
7	Thursday	1	0	3	11	14
8	Tuesday	1	0	4	10	14
9	Tuesday	1	1	1	7	8
10	Tuesday	0	1	0	9	9



# Travaux pratiques

Écrire le modèle qui prédit la présence en fonction du rappel.

$$y_i \sim \text{Binomial}(n_i, p_i)$$

$$\text{logit}(p_i) = \alpha + \beta \times \text{reminder}_i$$

$$\alpha \sim \text{Normal}(0, 1)$$

$$\beta \sim \text{Normal}(0, 1)$$



# Travaux pratiques

Écrire le modèle qui prédit la présence en fonction du rappel.

```
1 priors <- c(  
2   prior(normal(0, 1), class = Intercept),  
3   prior(normal(0, 1), class = b)  
4 )  
5  
6 mod8 <- brm(  
7   presence | trials(total) ~ 1 + reminder,  
8   family = binomial(link = "logit"),  
9   prior = priors,  
10  data = df3,  
11  cores = parallel::detectCores()  
12 )
```



# Travaux pratiques

Quel est l'effet **relatif** du mail de rappel ?

```
1 exp(fixef(mod8)[2]) # rapport des cotes sans vs. avec mail de rappel
```

```
[1] 2.986724
```

Envoyer un mail de rappel augmente la cote (le rapport des chances) par environ 3.



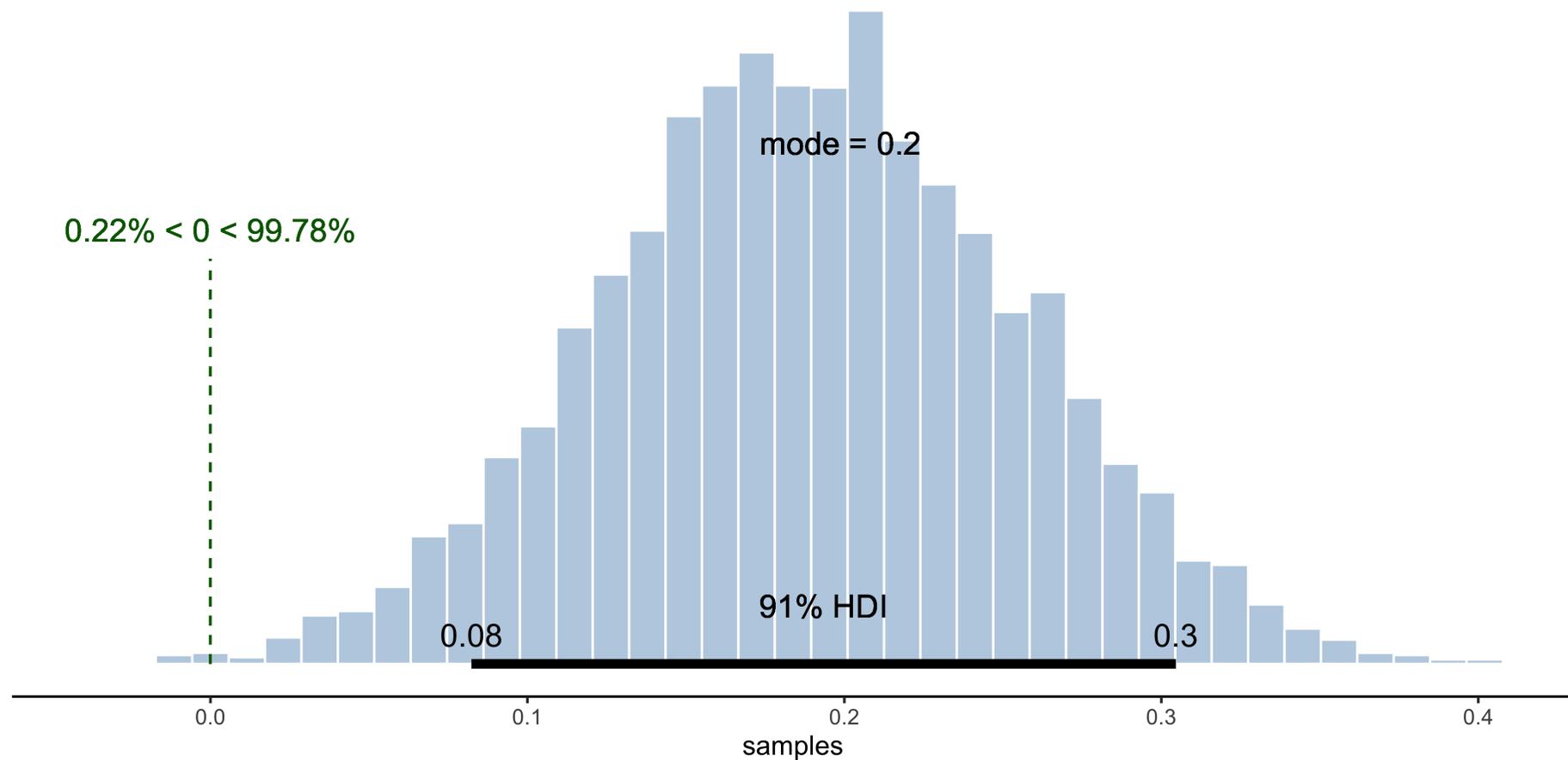
# Travaux pratiques

Quel est l'effet **absolu** du mail de rappel ?

```

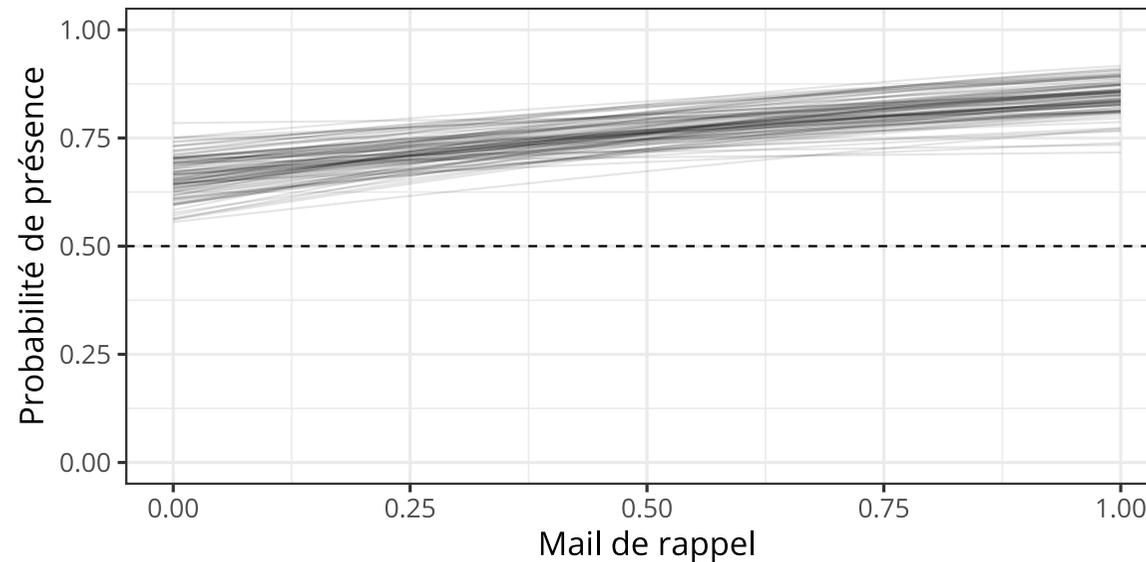
1 post <- as_draws_df(x = mod8) # récupères les échantillons du posterior
2 p.no <- plogis(post$b_Intercept) # probabilité de présence sans mail de rappel
3 p.yes <- plogis(post$b_Intercept + post$b_reminder) # probabilité de présence avec mail de rappel
4 posterior_plot(samples = p.yes - p.no, compval = 0, usemode = TRUE)

```



# Travaux pratiques

```
1 library(tidybayes)
2 library(modelr)
3
4 df3 %>%
5   group_by(total) %>%
6   data_grid(reminder = seq_range(reminder, n = 1e2) ) %>%
7   add_fitted_draws(mod8, newdata = ., n = 100, scale = "linear") %>%
8   mutate(estimate = plogis(.value) ) %>%
9   group_by(reminder, .draw) %>%
10  summarise(estimate = mean(estimate) ) %>%
11  ggplot(aes(x = reminder, y = estimate, group = .draw) ) +
12  geom_hline(yintercept = 0.5, lty = 2) +
13  geom_line(aes(y = estimate, group = .draw), size = 0.5, alpha = 0.1) +
14  ylim(0, 1) +
15  labs(x = "Mail de rappel", y = "Probabilité de présence")
```



# Travaux pratiques

- Quelle est la probabilité qu'un participant, qui s'est inscrit de son propre chef, vienne effectivement passer l'expérience ?
- Quel est l'effet du rappel ?
- **Quel est l'effet du mode d'inscription ?**
- Quel est l'effet conjoint de ces deux prédicteurs ?



# Travaux pratiques

Écrire le modèle qui prédit la présence en fonction du mode d'inscription.

$$y_i \sim \text{Binomial}(n_i, p_i)$$

$$\text{logit}(p_i) = \alpha + \beta \times \text{inscription}_i$$

$$\alpha \sim \text{Normal}(0, 1)$$

$$\beta \sim \text{Normal}(0, 1)$$



# Travaux pratiques

```
1 priors <- c(  
2   prior(normal(0, 1), class = Intercept),  
3   prior(normal(0, 1), class = b)  
4 )  
5  
6 mod9 <- brm(  
7   presence | trials(total) ~ 1 + inscription,  
8   family = binomial(link = "logit"),  
9   prior = priors,  
10  data = df3,  
11  cores = parallel::detectCores()  
12 )
```

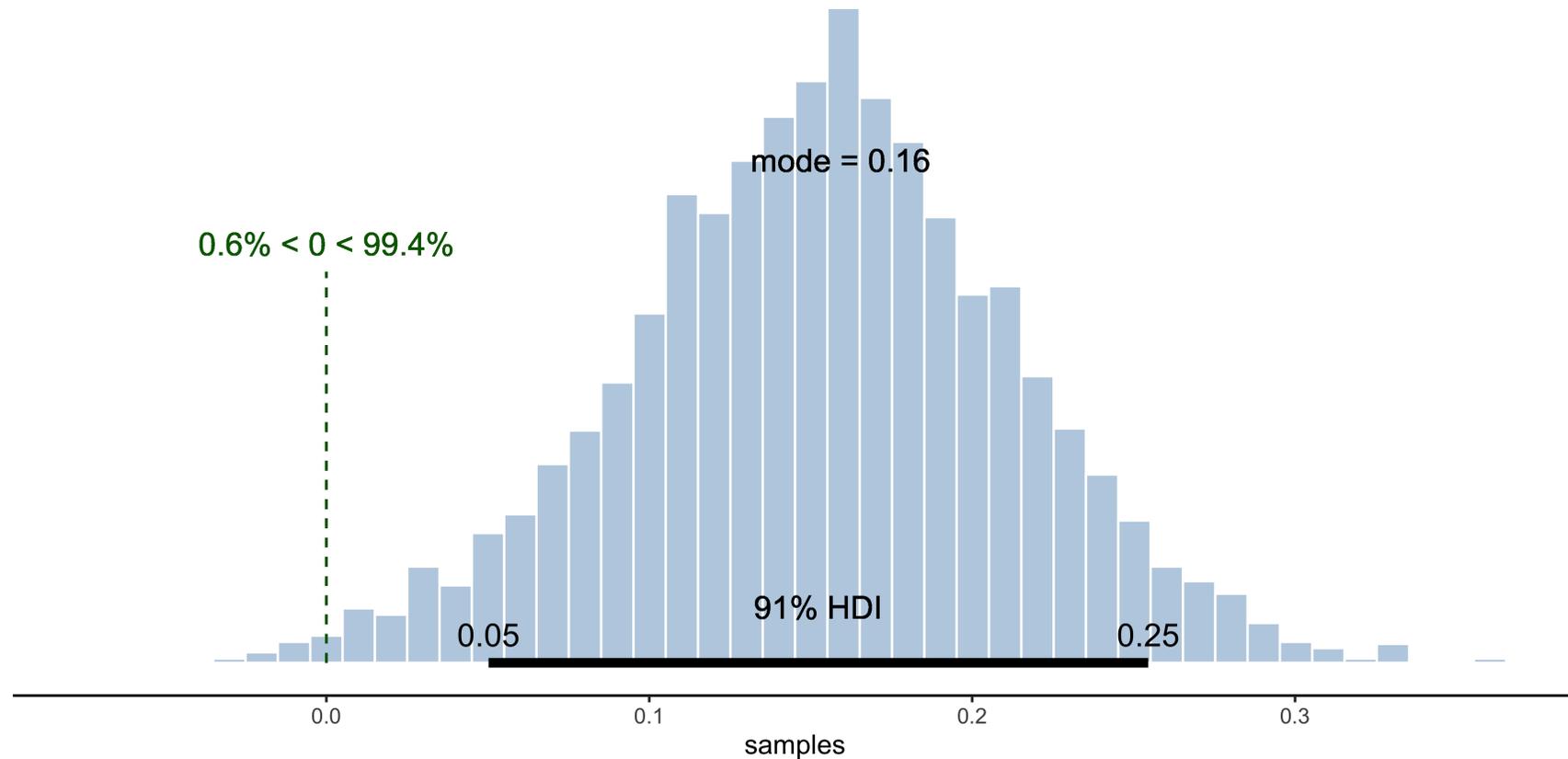


# Travaux pratiques

```

1 post <- as_draws_df(x = mod9)
2 p.panel <- plogis(post$b_Intercept) # probabilité moyenne de présence - panel
3 p.doodle <- plogis(post$b_Intercept + post$b_inscription) # probabilité moyenne de présence - doodle
4 posterior_plot(samples = p.panel - p.doodle, compval = 0, usemode = TRUE)

```



La probabilité de présence est augmentée d'environ 0.16 lorsque l'on s'inscrit sur un panel comparativement à une inscription sur un Doodle (effet légèrement plus faible que pour le rappel).



# Travaux pratiques

- Quelle est la probabilité qu'un participant, qui s'est inscrit de son propre chef, vienne effectivement passer l'expérience ?
- Quel est l'effet du rappel ?
- Quel est l'effet du mode d'inscription ?
- **Quel est l'effet conjoint de ces deux prédicteurs ?**



# Travaux pratiques

Écrire le modèle complet.

$$y_i \sim \text{Binomial}(n_i, p_i)$$

$$\text{logit}(p_i) = \alpha + \beta_1 \times \text{reminder}_i + \beta_2 \times \text{inscription}_i$$

$$\alpha \sim \text{Normal}(0, 1)$$

$$\beta_1, \beta_2 \sim \text{Normal}(0, 1)$$



# Travaux pratiques

```
1 priors <- c(  
2   prior(normal(0, 1), class = Intercept),  
3   prior(normal(0, 1), class = b)  
4 )  
5  
6 mod10 <- brm(  
7   presence | trials(total) ~ 1 + reminder + inscription,  
8   family = binomial(link = "logit"),  
9   prior = priors,  
10  data = df3,  
11  cores = parallel::detectCores()  
12 )
```



# Travaux pratiques

```
1 summary(mod10)
```



```
Family: binomial
Links: mu = logit
Formula: presence | trials(total) ~ 1 + reminder + inscription
Data: df3 (Number of observations: 13)
Draws: 4 chains, each with iter = 2000; warmup = 1000; thin = 1;
       total post-warmup draws = 4000

Population-Level Effects:
      Estimate Est.Error 1-95% CI u-95% CI Rhat Bulk_ESS Tail_ESS
Intercept      1.01      0.57  -0.11   2.10 1.00    2311    2244
reminder       0.91      0.49  -0.00   1.90 1.00    2140    2111
inscription    -0.34      0.53  -1.37   0.71 1.00    2291    2353

Draws were sampled using sampling(NUTS). For each parameter, Bulk_ESS
and Tail_ESS are effective sample size measures, and Rhat is the potential
scale reduction factor on split chains (at convergence, Rhat = 1).
```



# Travaux pratiques

Le mail de rappel semble avoir moins d'effet dans le modèle complet que dans le modèle simple... pourquoi ?

```
1 fixef(mod8) %>% exp() # calcul du "odds ratio" 
```

	Estimate	Est.Error	Q2.5	Q97.5
Intercept	1.964499	1.274686	1.236716	3.188580
reminder	2.986724	1.478333	1.404487	6.698153

```
1 fixef(mod9) %>% exp() # calcul du "odds ratio" 
```

	Estimate	Est.Error	Q2.5	Q97.5
Intercept	6.221081	1.434904	3.1658702	12.8257873
inscription	0.384240	1.501875	0.1682796	0.8497573

```
1 fixef(mod10) %>% exp() # calcul du "odds ratio" 
```

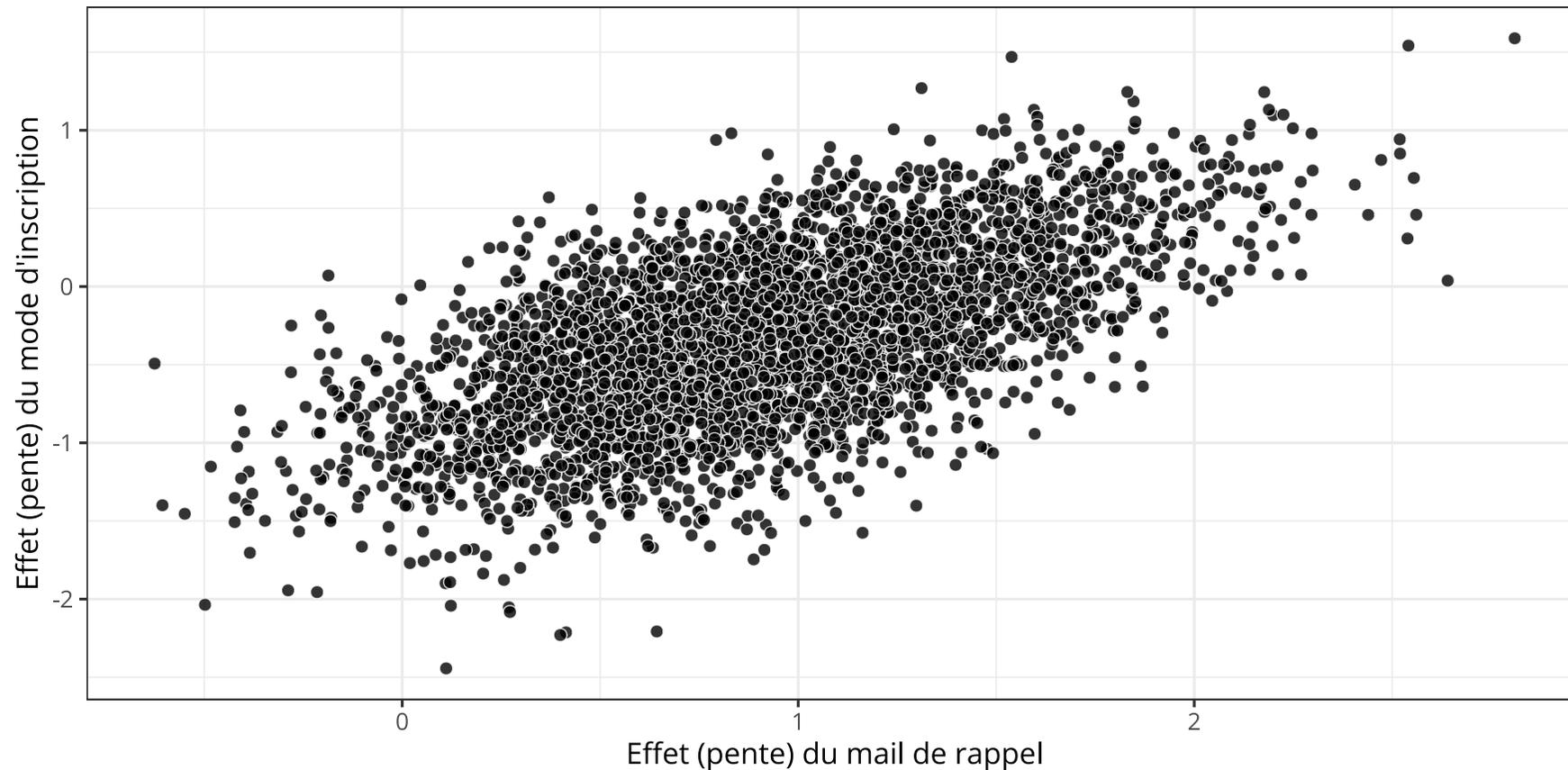
	Estimate	Est.Error	Q2.5	Q97.5
Intercept	2.7442123	1.769934	0.8921157	8.166210
reminder	2.4807285	1.626489	0.9973190	6.698664
inscription	0.7102085	1.707242	0.2538193	2.023864



# Travaux pratiques

On a déjà rencontré ce cas de figure (cf. Cours n°04). Lorsque deux prédicteurs contiennent une part d'information commune, l'estimation des pentes est corrélée...

```
1 as_draws_df(x = mod10) %>%  
2   ggplot(aes(b_reminder, b_inscription) ) +  
3   geom_point(size = 3, pch = 21, alpha = 0.8, color = "white", fill = "black") +  
4   labs(x = "Effet (pente) du mail de rappel", y = "Effet (pente) du mode d'inscription")
```



# Travaux pratiques

En effet, les données ont été collectées par deux expérimentateurs. L'un d'entre eux a recruté tous ses participants via Doodle, et n'envoyait pas souvent de mail de rappel. Le deuxième expérimentateur a recruté tous ses participants via un panneau physique présent dans le laboratoire et envoyait systématiquement un mail de rappel. Autrement dit, ces deux variables sont presque parfaitement confondues.

```
1 open_data(absence) %>%
2   group_by(inscription, reminder) %>%
3   summarise(n = sum(total) ) %>%
4   spread(key = reminder, value = n) %>%
5   data.frame()
```

```
inscription no yes
1      doodle 72  22
2      panel NA  51
```

