Bayesian data analysis: Theory & practice

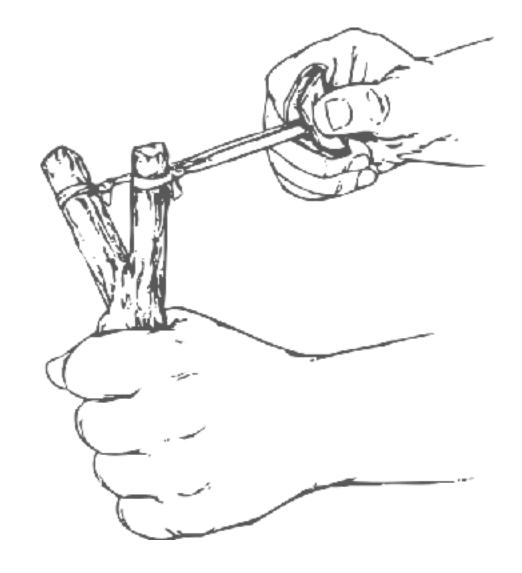
Part 3a: Categorical predictors & generalized linear models

Michael Franke

Main learning goals

for this part

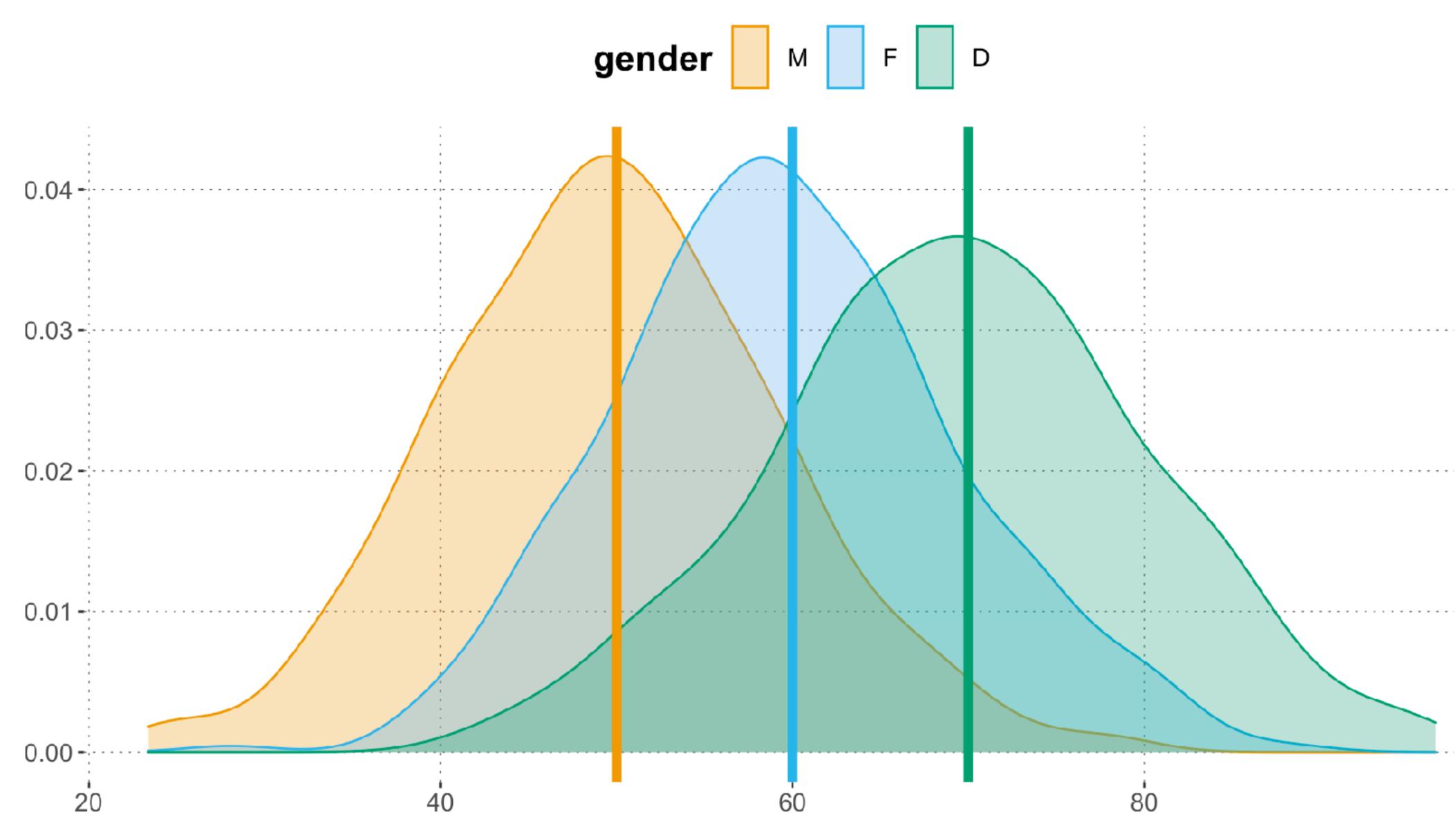
- 1. multiple regression w/ categorical predictors
 - a. contrast coding
 - b. derived variables in Bayesian analysis
- 2. generalized linear models
 - a. logistic
 - b. multinomial
 - c. ordinal
- 3. excursion: "beyond GLMs"
 - a. mixture models
 - b. distributional models



contrast coding

Some fake data

three-way categorical variable



Finding numbers for categories

metric predictors

У	Χo	X 1	X ₂
42	1	4	163
19	1	7	128
38	1	2	99
•	•	•	•

categorical predictor

У	gender
51	M
59	F
73	D
•	•

Treatment coding

comparing against a reference category

У	gender	Xo		
51	M	1	0	0
59	F		1	0
73	D		0	1
•	•	•		•

$$\frac{2}{2}$$

$$\hat{\mu}_{M} = \beta_{0}1 + \beta_{1}0 + \beta_{2}0$$

$$\hat{\mu}_{F} = \beta_{0}1 + \beta_{1}1 + \beta_{2}0$$

$$\hat{\mu}_{D} = \beta_{0}1 + \beta_{1}0 + \beta_{2}1$$

$$y_i = \sum_{j=0}^{k} \beta_j x_{ij} + \epsilon_i$$

hypotheses

$$\beta_0 = \hat{\mu}_M$$

$$\beta_1 = \hat{\mu}_F - \hat{\mu}_M$$

$$\beta_2 = \hat{\mu}_D - \hat{\mu}_M$$

A STATE OF STATE OF THE STATE O

Treatment coding

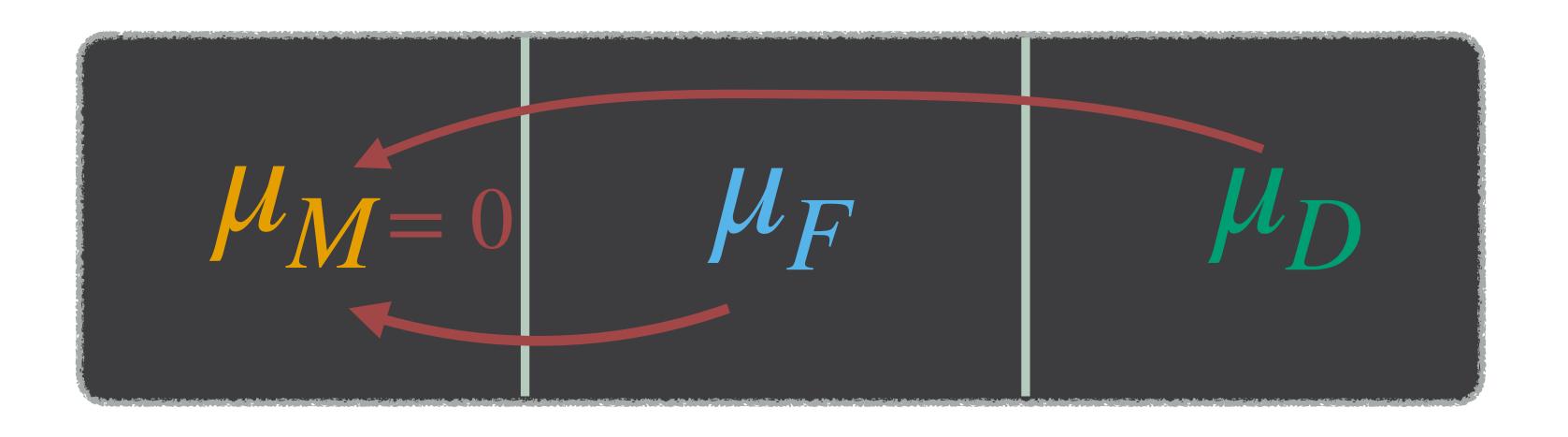
comparing against a reference category

У	gender	Χo	X 1	X ₂
51	M	1	0	0
59	F	1	1	0
73	D	1	0	1
•	•	•	•	•

$$\beta_0 = \hat{\mu}_M$$

$$\beta_1 = \hat{\mu}_F - \hat{\mu}_M$$

$$\beta_2 = \hat{\mu}_D - \hat{\mu}_M$$



Cell means coding

estimating a mean for each cell

У	gender	Χo	X ₁	X ₂
51	M	1	0	0
59	F	0	1	0
73	D	0	0	1
•	•	•	•	•

$$\beta_0 = \hat{\mu}_M$$

$$\beta_1 = \hat{\mu}_F$$

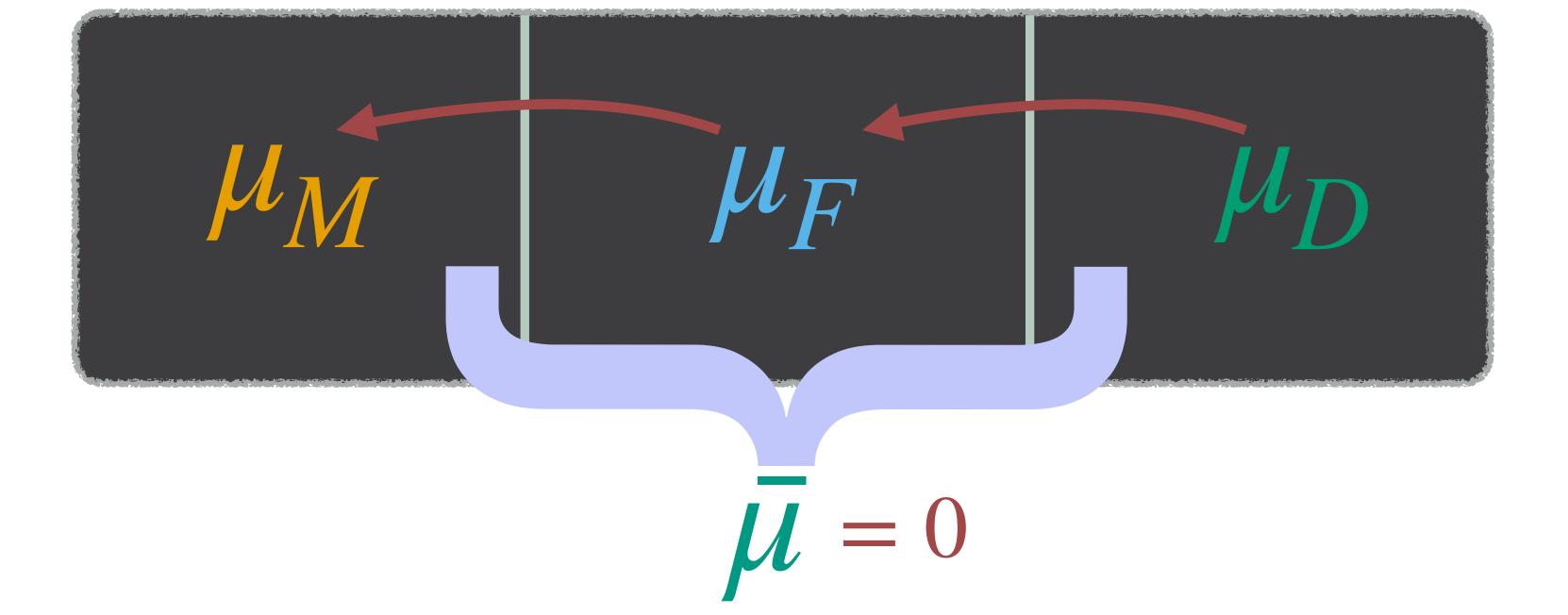
$$\beta_2 = \hat{\mu}_D$$

$$\mu_{M}=0 \qquad \mu_{F}=0 \qquad \mu_{D}=0$$

Simple difference coding

estimating a mean for each cell

У	gender	Χo	X ₁	X ₂
51	M	1	-2/3	-1/ ₃
59	F	1	1/3	_1/3
73	D	1	1/3	2/3
•	•	•	•	•



$$\beta_0 = \bar{\mu}$$

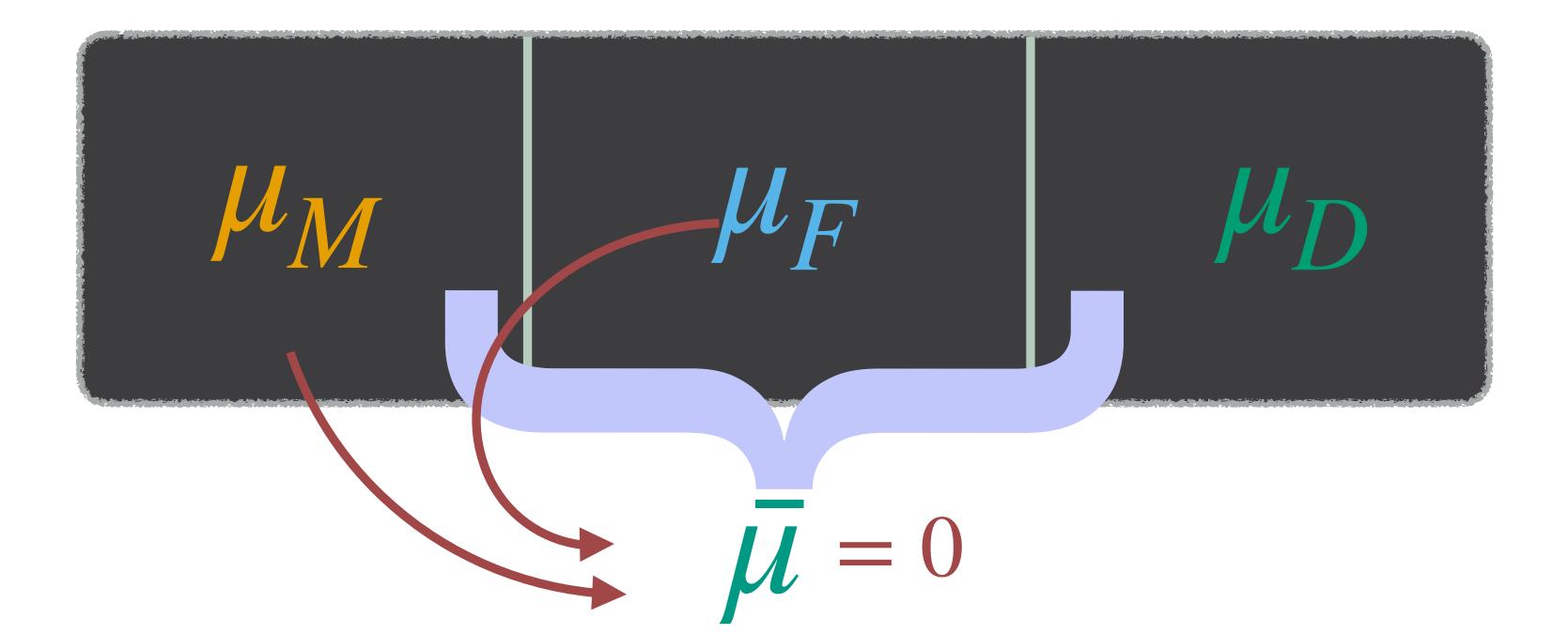
$$\beta_1 = \hat{\mu}_F - \hat{\mu}_M$$

$$\beta_2 = \hat{\mu}_D - \hat{\mu}_F$$

Sum coding

comparing k-1 cells to the grand mean

У	gender	Χo	X 1	X ₂
51	M	1	1	0
59	F	1	0	1
73	D	1	-1	-1
•	•	•	•	•



$$\beta_0 = \bar{\mu}$$

$$\beta_1 = \hat{\mu}_M - \bar{\mu}$$

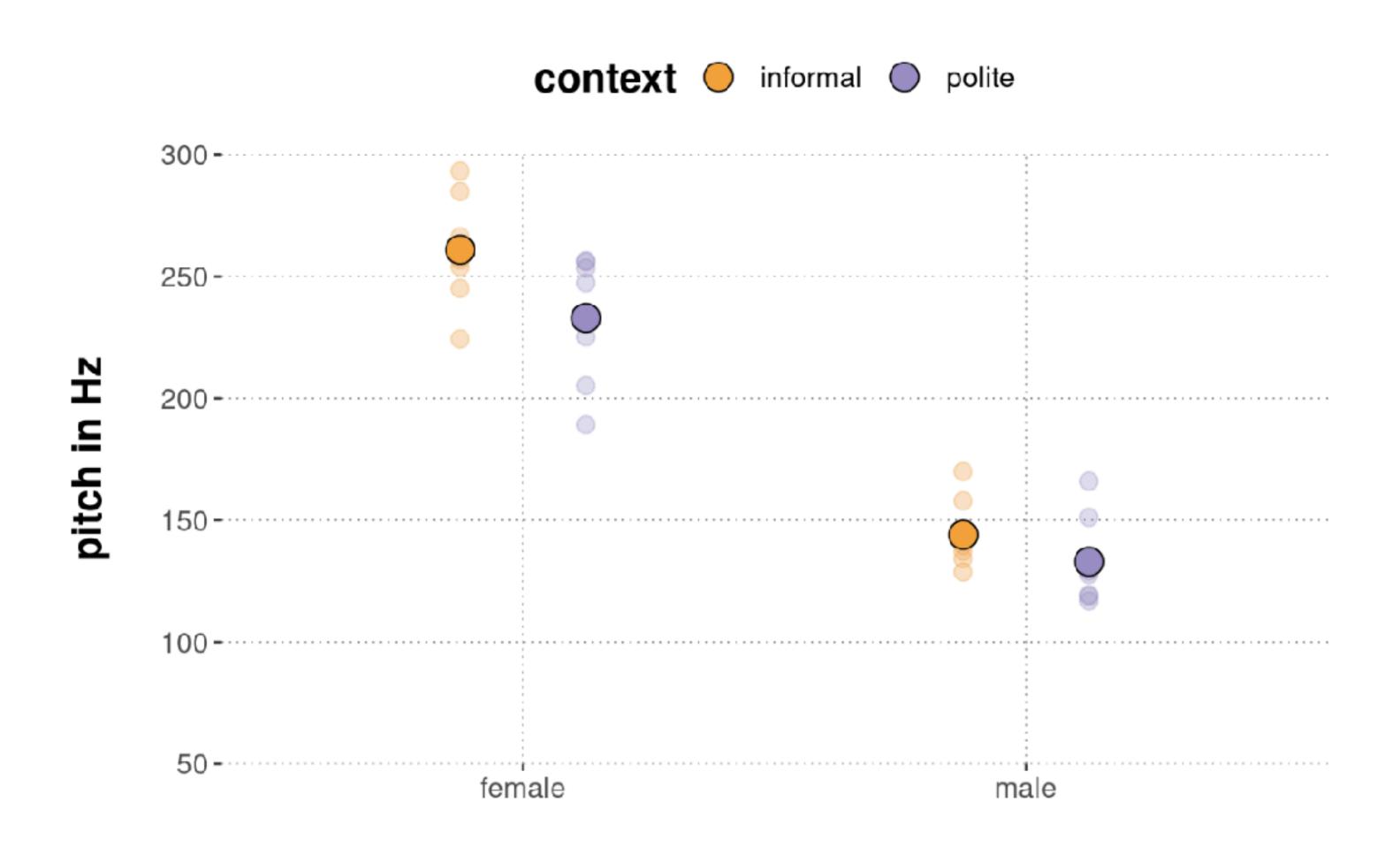
$$\beta_2 = \hat{\mu}_F - \bar{\mu}$$

Case study: pitch in context

data from Winter & Grawunder (2012)

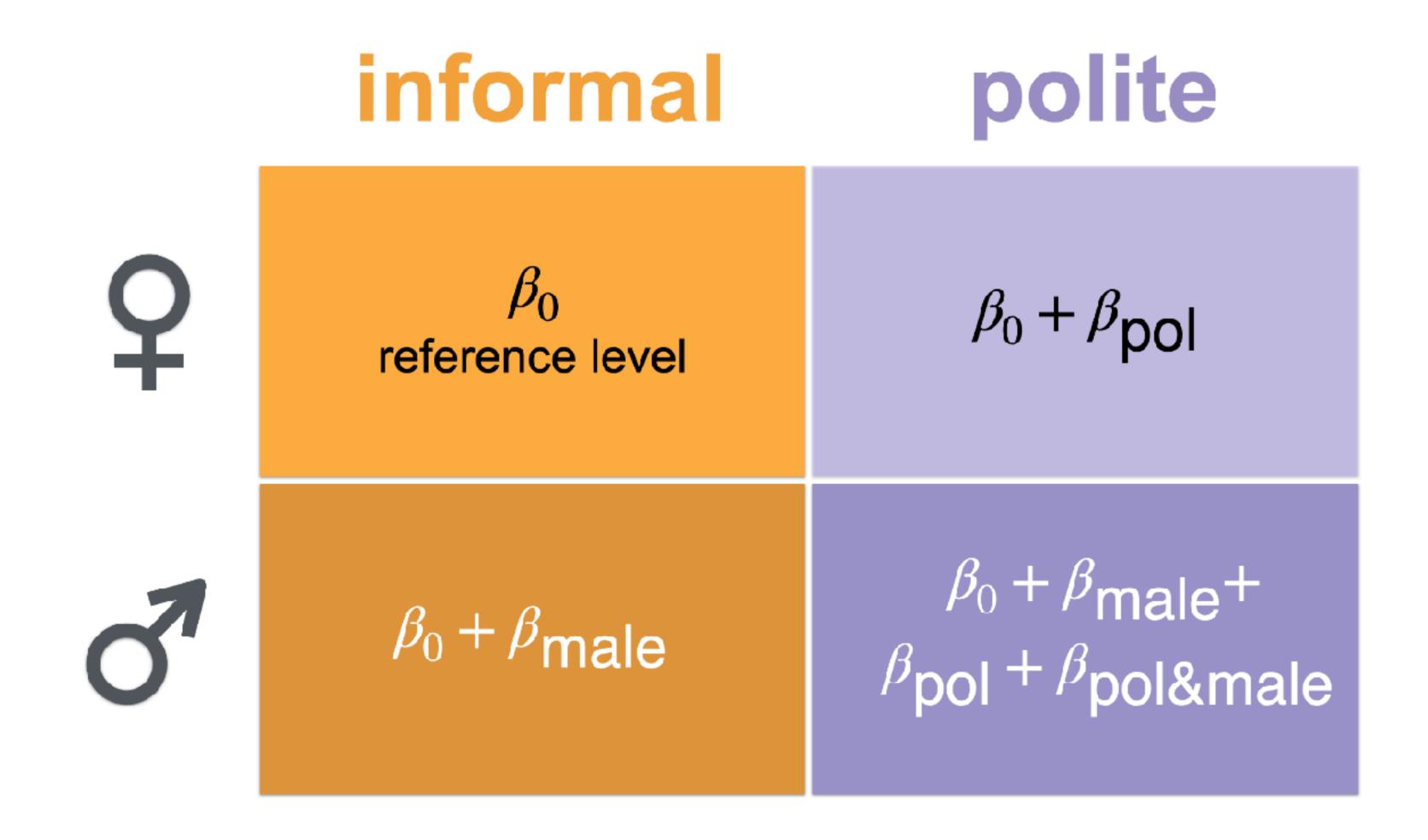
```
politeness_data <- aida::data_polite
politeness_data %>% head(5)
```

```
## # A tibble: 5 × 5
     subject gender sentence context pitch
    <chr> <chr> <chr>
                            <chr>
                                     <dbl>
## 1 F1
                                     213.
                             pol
## 2 F1
                            inf
                                     204.
                                      285.
                             pol
                                      260.
## 5 F1
                                      204.
                             pol
```



read more <u>here</u>

Dummy coding for 2x2 design

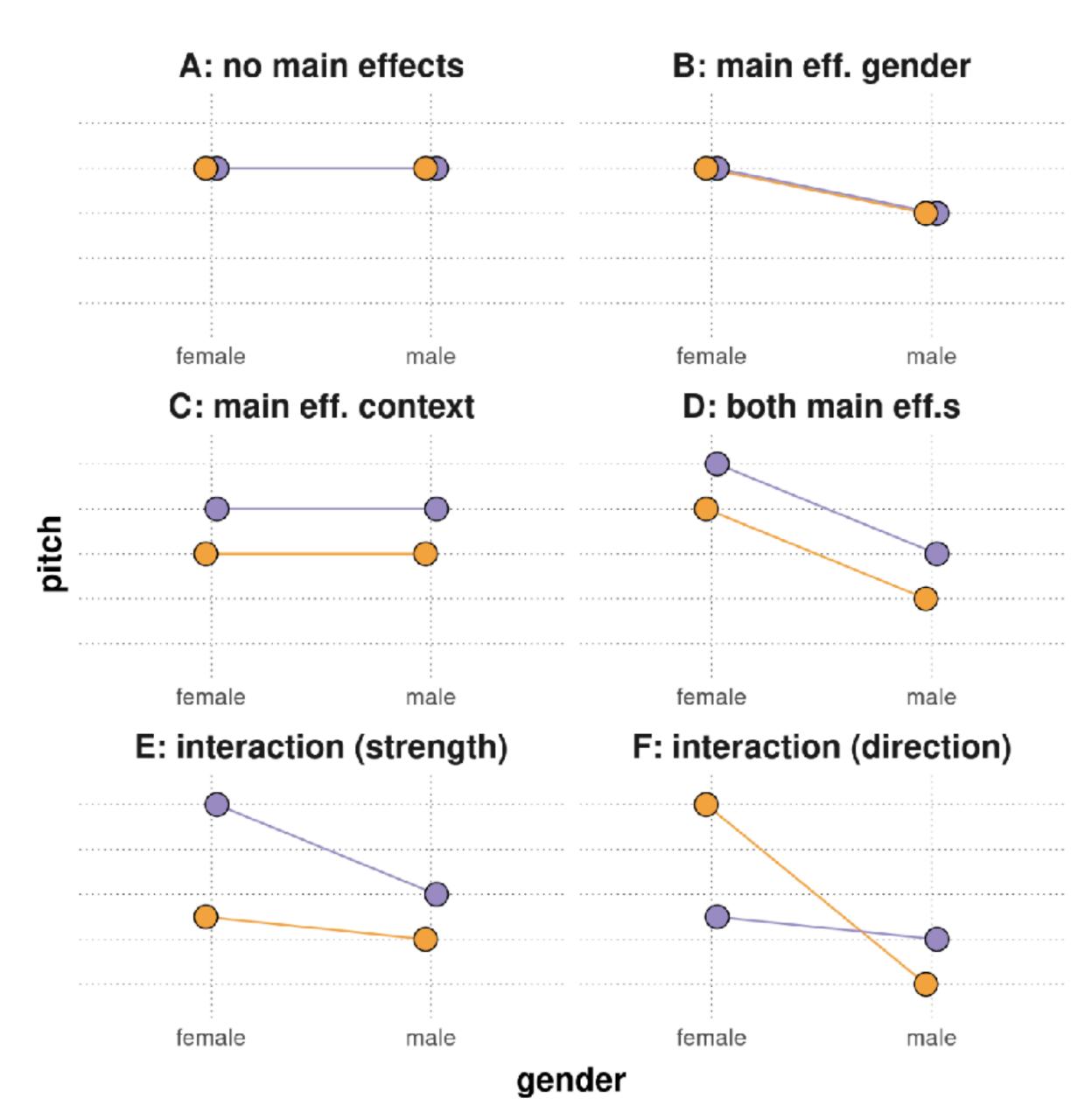


Interaction term $\beta_{\text{pol\&male}}$ is a 'difference of differences'

Main effects & interactions

in 2x2 designs

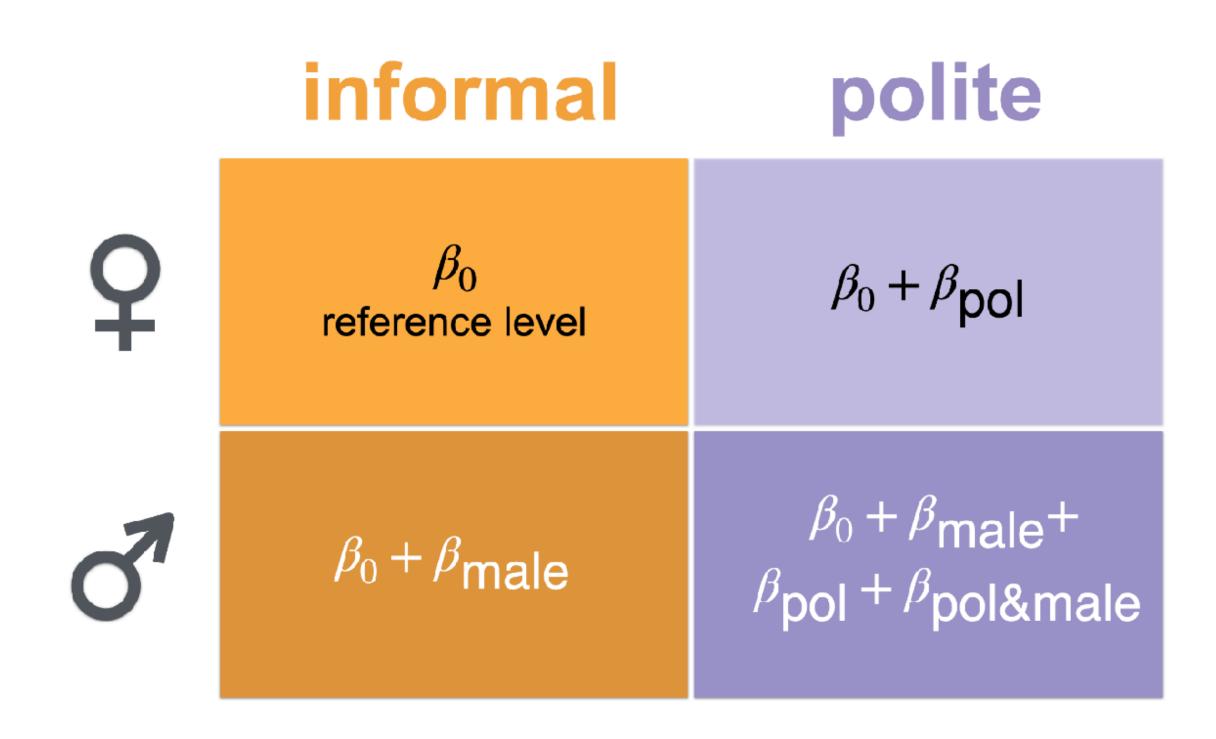
Interaction term $\beta_{\rm pol\&male}$ is a 'difference of differences'



context informal polite

Bayesian regression

```
# here, we only use fixed effects
fit_dummy_FE <- brm(
  pitch ~ gender * context,
  data = politeness_df,
  cores = 4,
  iter = 1000
)</pre>
```



Population-Level Effects:

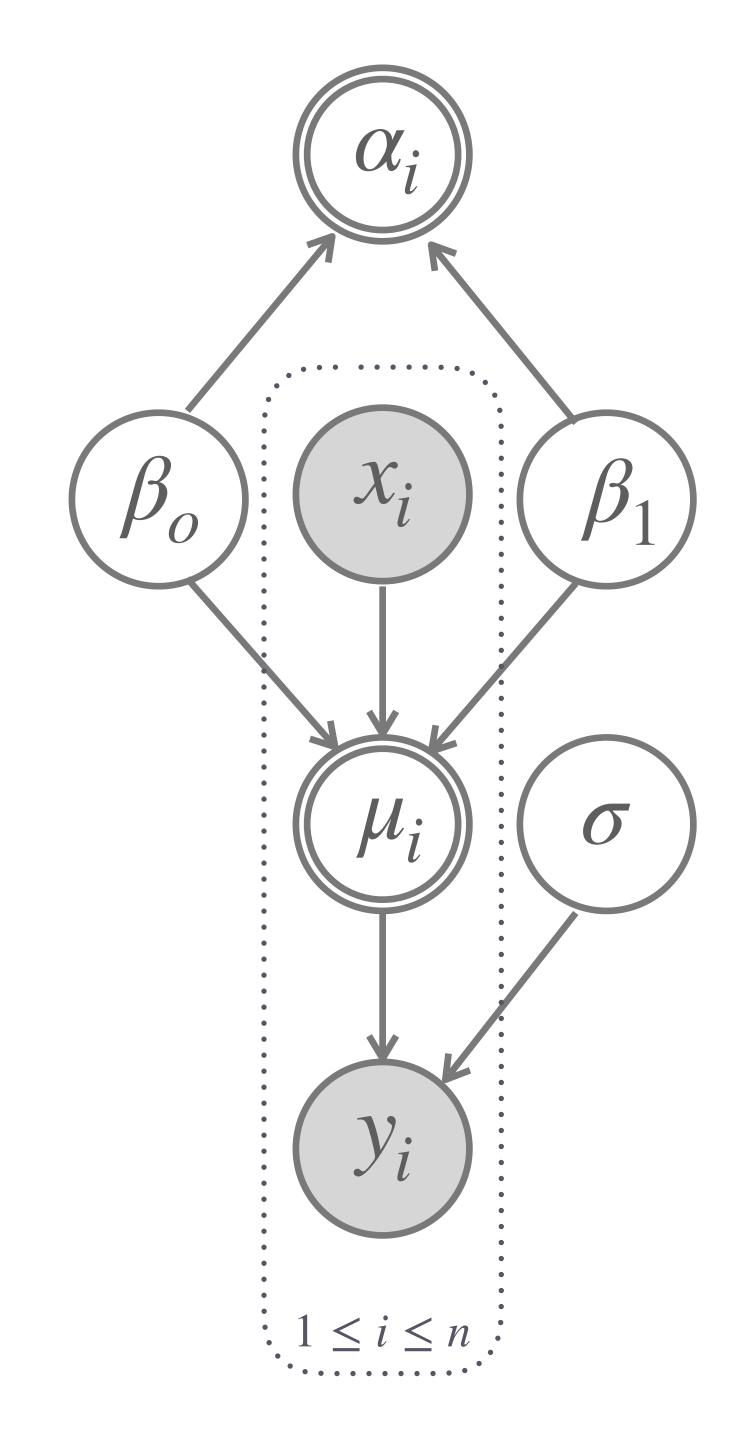
	Estimate	Est.Error	l-95% CI	u−95% CI
Intercept	260.56	7.87	244.15	275.21
genderM	-116.16	11.01	-137.31	-94.05
contextpol	-27.23	11.10	-48.38	-5.23
<pre>genderM:contextpol</pre>	15.77	16.05	-16.54	46.24

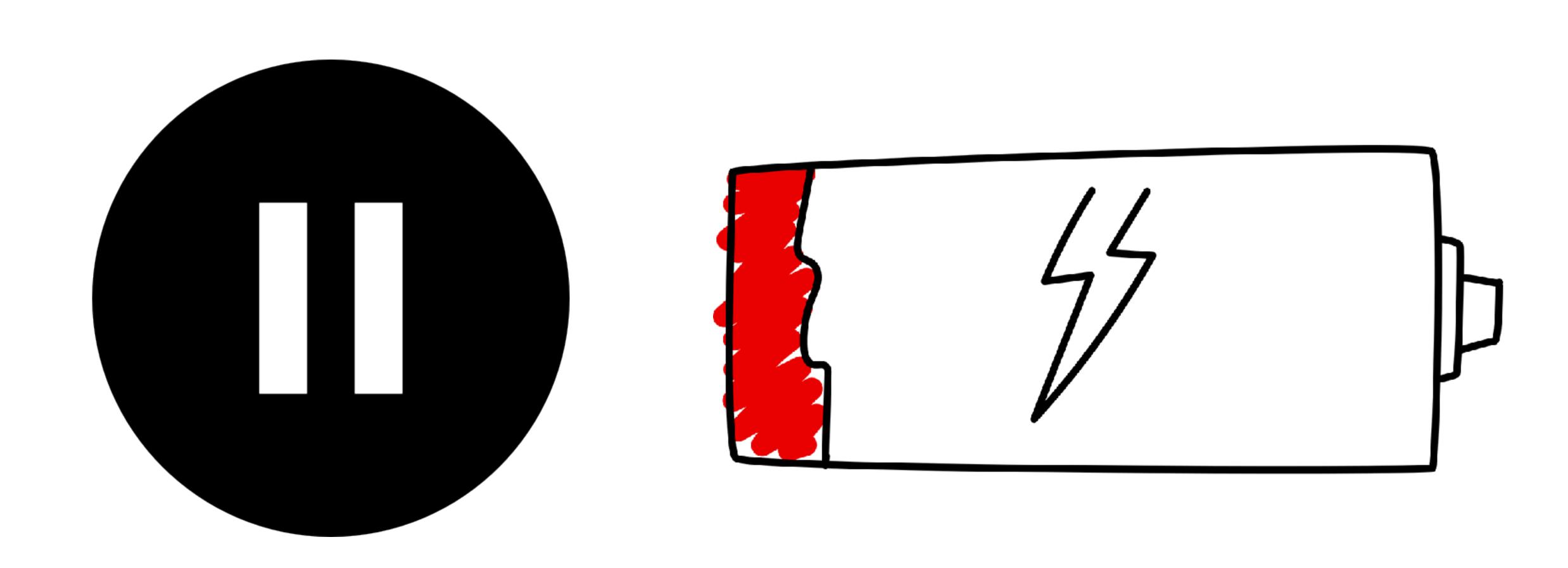
Which questions about cell mean differences can we address with this information directly?



Derived variables

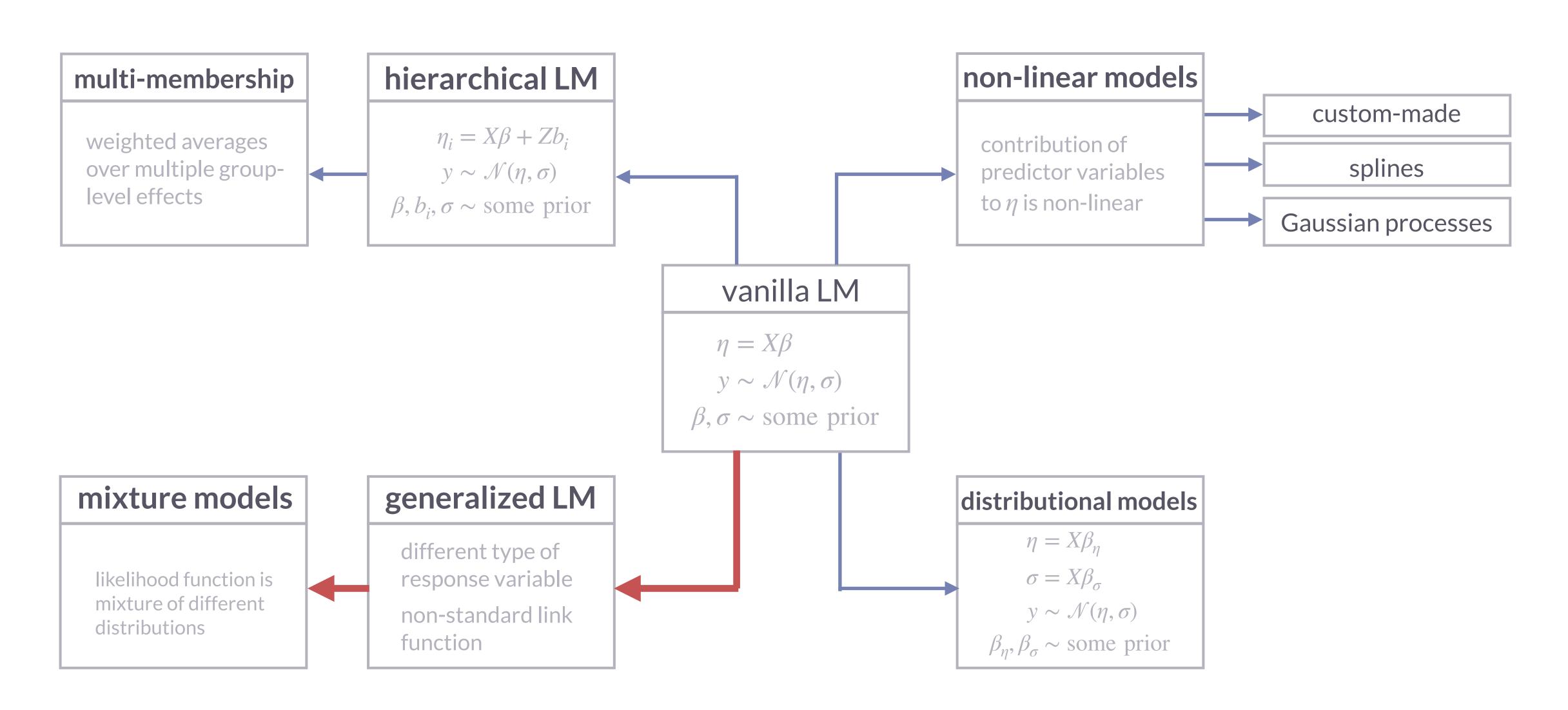
- obtain samples from model parameters
- apply (deterministic) function to each sample
 - to derive (deterministically) a new model variable
- violà: samples from the posterior of a new "derived variable"





Roadmap "beyond vanilla"

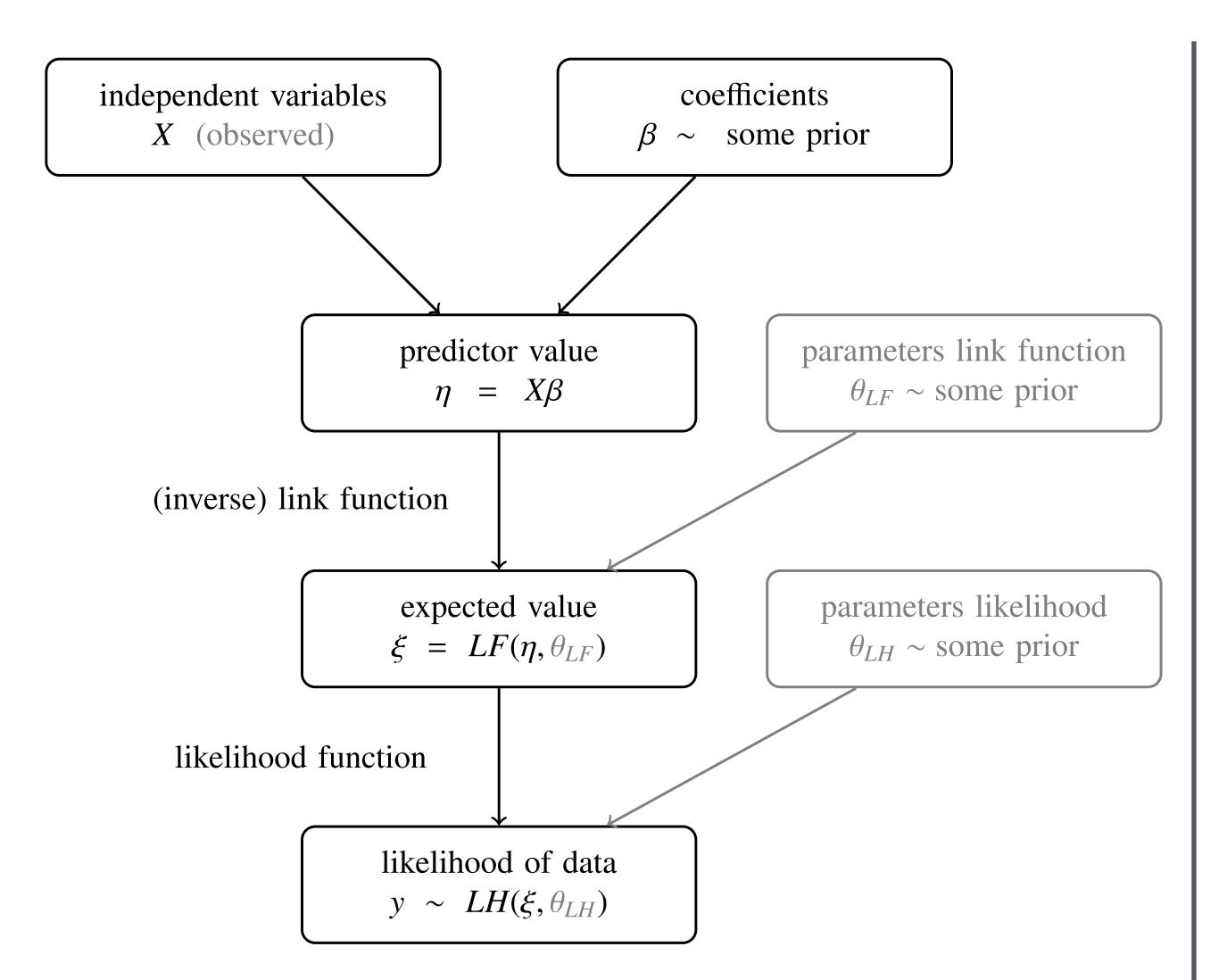
common extensions of linear regression modeling



Generalized linear regression models

Generalized linear regression model

Coefficients → linear predictor → central tendency → likelihood



Simple linear regression

$$\eta_i = \mathbf{x}_i \cdot \boldsymbol{\beta}$$
 [linear predictor]

$$\xi_i = \eta_i$$
 [predictor of central tendency]

$$y_i \sim \text{Normal}(\xi_i, \sigma)$$
 [likelihood]

Logistic regression

$$\eta_i = \mathbf{x}_i \cdot \boldsymbol{\beta}$$
 [linear predictor]

$$\xi_i = \text{logistic}(\eta_i)$$
 [predictor of central tendency]

$$y_i \sim \text{Bernoulli}(\xi_i)$$
 [likelihood]

Poisson regression

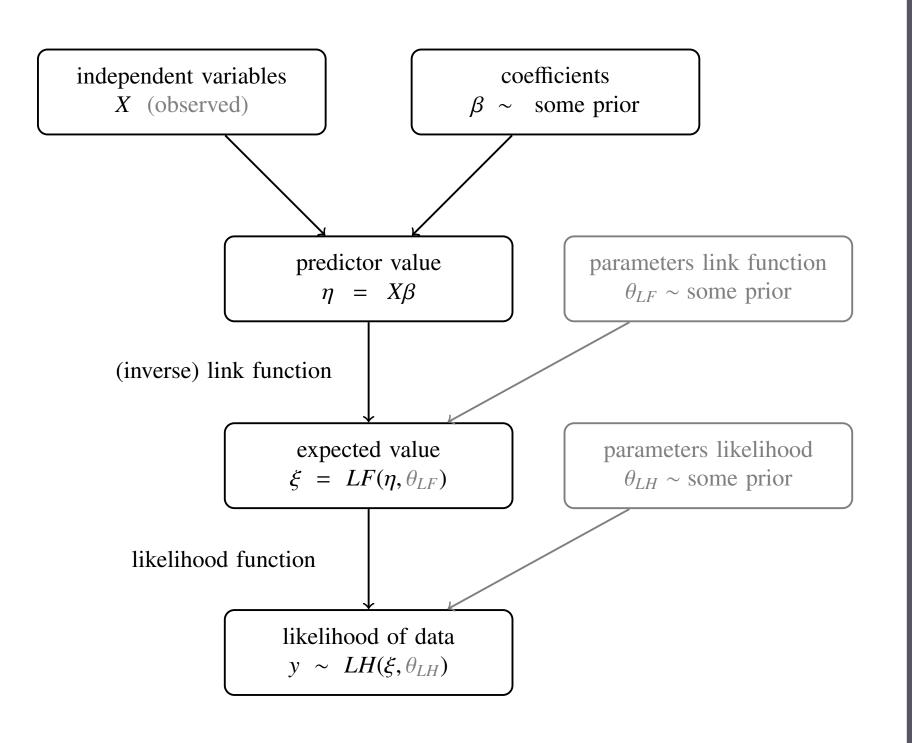
$$\eta_i = \mathbf{x}_i \cdot \boldsymbol{\beta}$$
 [linear predictor]

$$\xi_i = \exp(\eta_i)$$
 [predictor of central tendency]

$$\xi_i = \exp(\eta_i)$$
 [predictor of $y_i \sim \operatorname{Poisson}(\xi_i)$ [likelihood]

Generalized linear regression model

Coefficients → linear predictor → central tendency → likelihood



$\mathbf{type}\ \mathbf{of}\ y$	(inverse) link function	likelihood function
metric	$\xi=\eta$	$y \sim \operatorname{Normal}(\xi; \sigma)$
binary	$\xi = \operatorname{logistic}(\eta)$	$y \sim \mathrm{Bernoulli}(\xi)$
nominal	$\xi = \operatorname{soft-max}(\eta)$	$y \sim \mathrm{Categorical}(\xi)$
ordinal	$\xi = \operatorname{cumulative-logit}(\eta; \delta)$	$y \sim \mathrm{Categorical}(\xi)$
count	$\xi=\exp(\eta)$	$y \sim \mathrm{Poisson}(\xi)$

read more <u>here</u>

The BRMS "family of families"

- link- and likelihood function are set by family parameter in brm function
- requires an object of type `brms:: brmsfamily`
 - many predefined families, listed <u>here</u>
 - instantiated by function calls like cumulative()
 - allow flexible parameterization
 - documentation of parameterization: <u>here</u>
 - creating custom families is possible, see here

```
fit_ordinal <- brm(
  formula = prototype_label ~ MAD,
  data = data_MT_prepped2,
  family = cumulative()
)</pre>
```

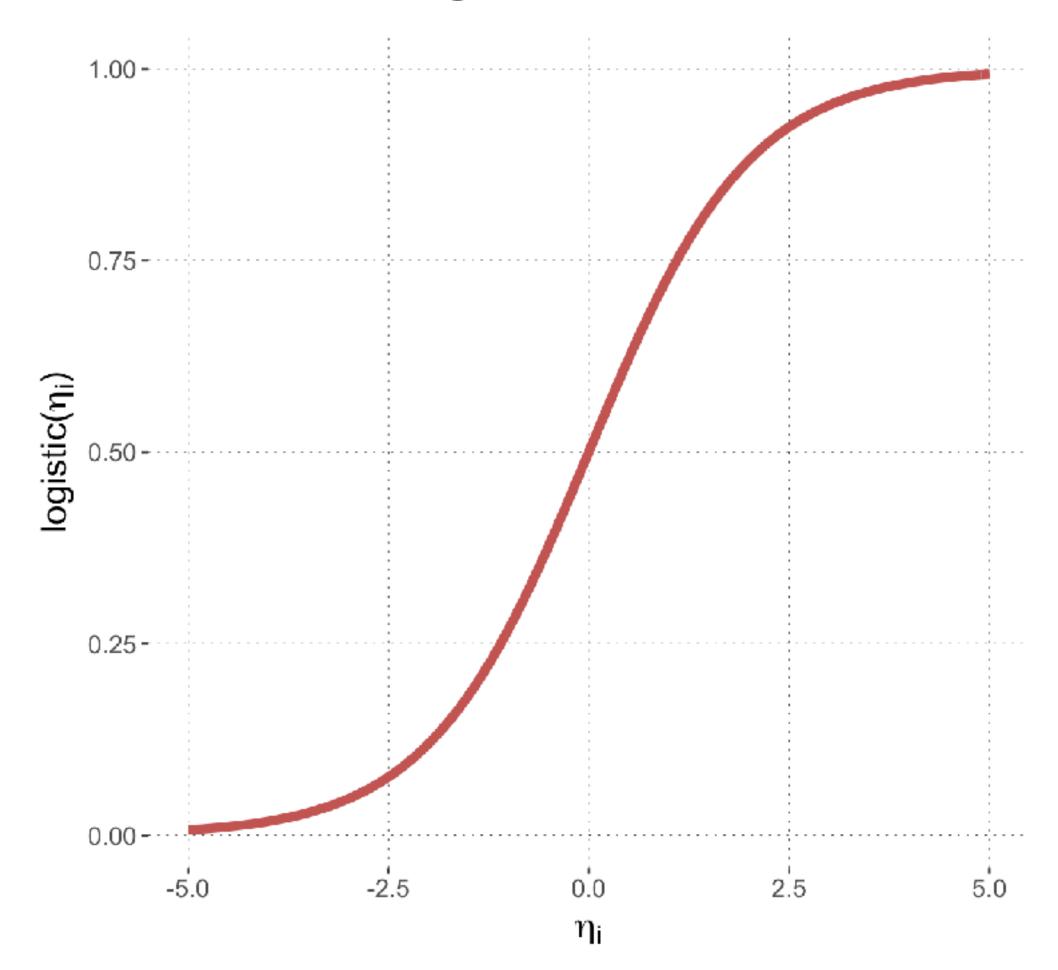
Logistic regression

Logistic regression

Definition

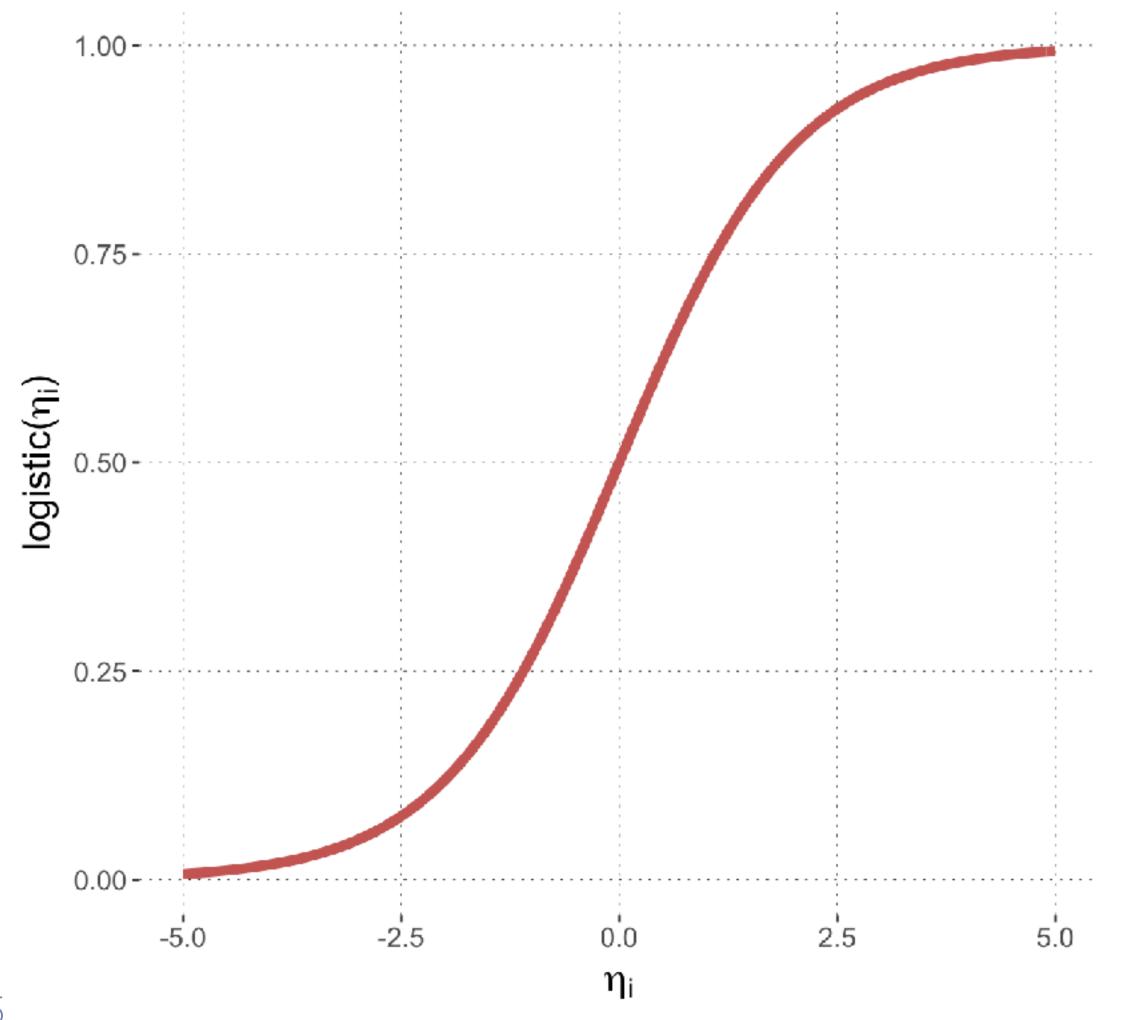
```
\begin{split} \eta_i &= \mathbf{x}_i \cdot \boldsymbol{\beta} & \text{[linear predictor]} \\ \xi_i &= \text{logistic}(\eta_i) & \text{[predictor of central tendency]} \\ y_i &\sim \text{Bernoulli}(\xi_i) & \text{[likelihood]} \end{split}
```

Link function: logistic

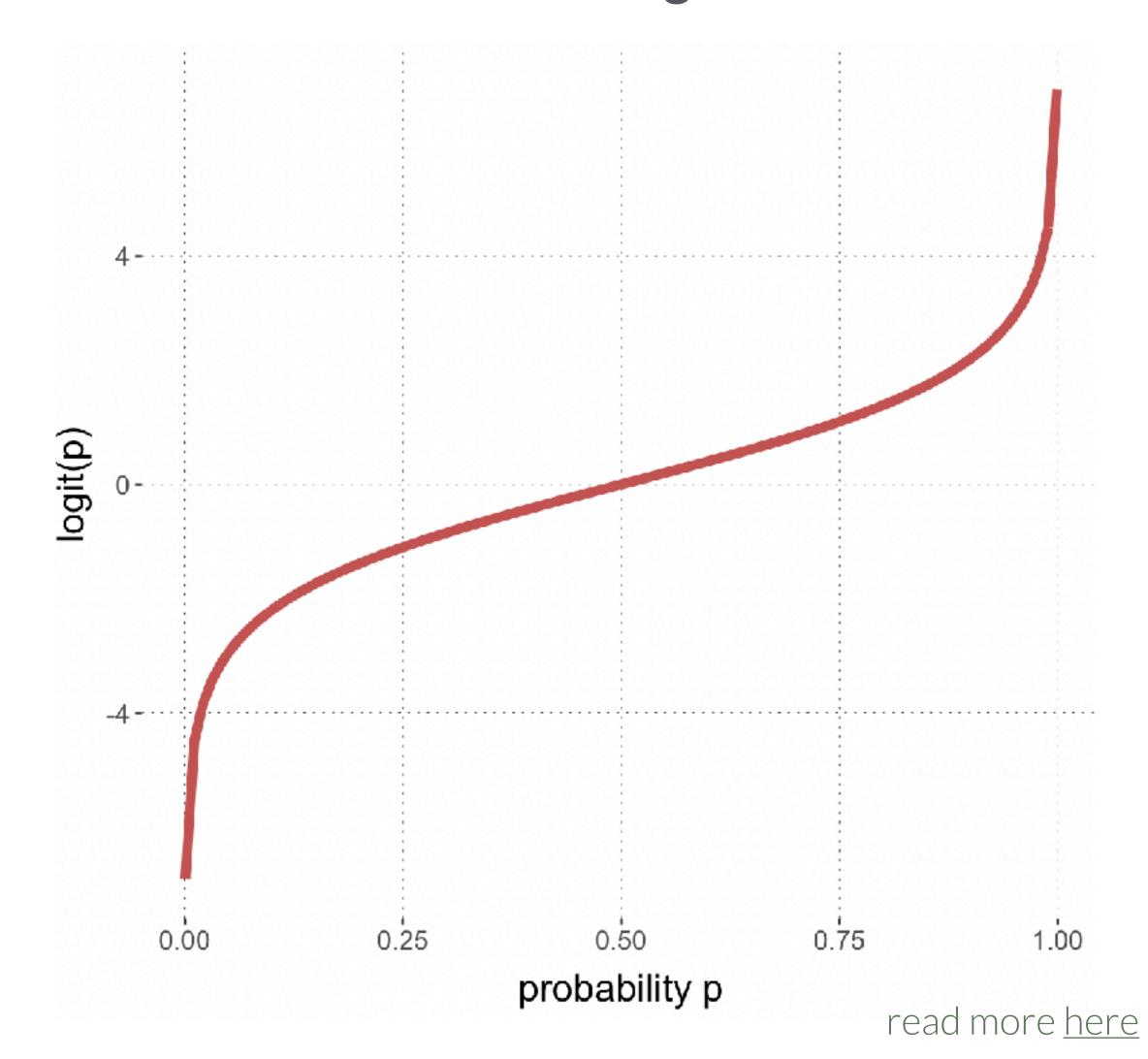


Logistic regression Link & inverse link function

Link function: logistic



"Inverse" link function: logit



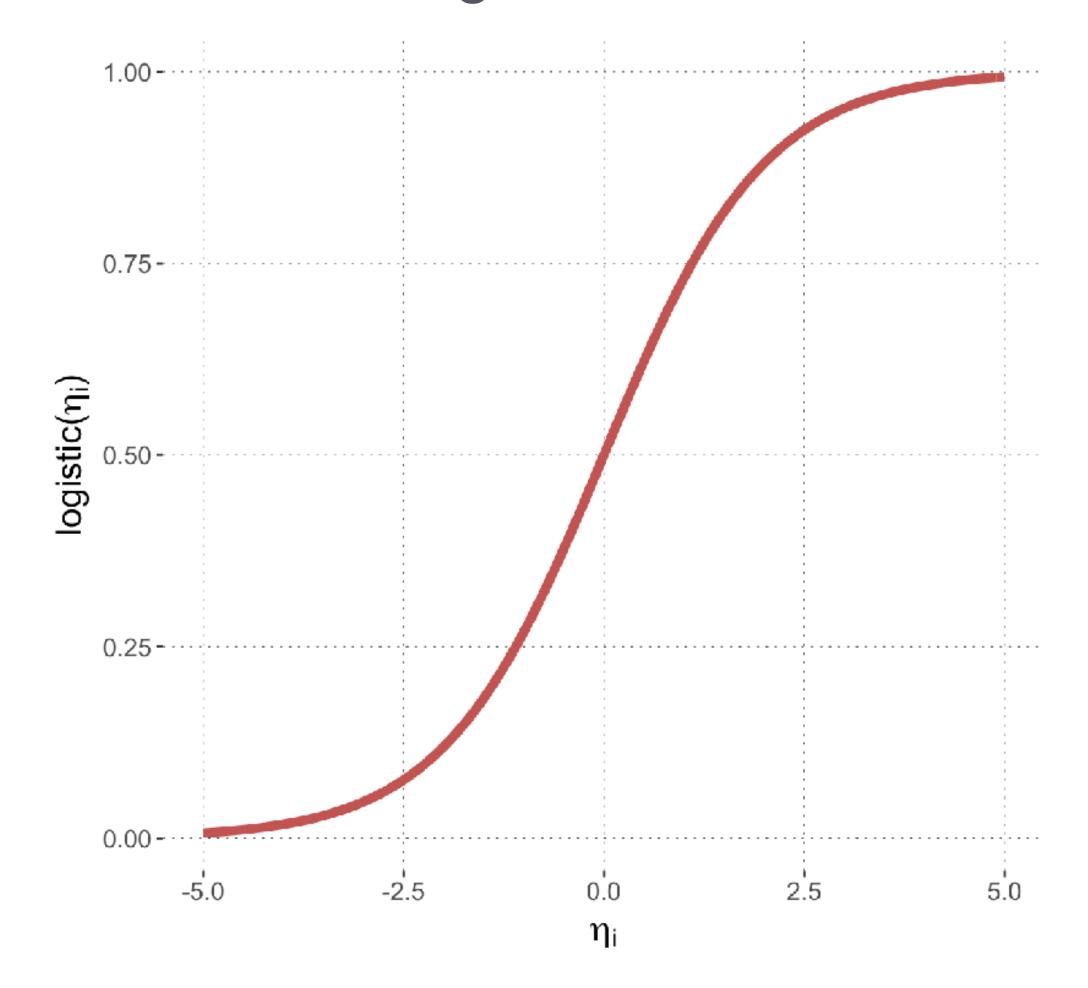
Logistic regression

interpretation

- linear predictor η encodes the log-odds
 - probability of "heads" vs "tails" (1 vs 0)
- unit change of η can be seen as a change in beliefs corresponding to a Bayes factor of ~2.72

$$\eta_1 - \eta_2 = \log \frac{\xi_1}{1 - \xi_1} - \log \frac{\xi_2}{1 - \xi_2} = \log \left(\frac{\xi_1}{1 - \xi_1} \frac{1 - \xi_2}{\xi_2} \right)
\Leftrightarrow \frac{\xi_1}{1 - \xi_1} = \exp(\eta_1 - \eta_2) \frac{\xi_2}{1 - \xi_2}$$

Link function: logistic





Multinomial regression

Multinomial regression

- we want to predict probabilities $\mathbf{p} = \langle p_1, ..., p_k \rangle$
 - ullet p_i is the prediction for category j probability of
- it suffices to estimate k-1 probabilities
 - probabilities sum to one
 - fix a reference category (similar to treatment coding!)
- (non-normalized) weights s_j from linear predictors:

$$s_i = \mathbf{x}_i \cdot \beta^j$$

probabilities from soft-max:

$$p_{j} = \frac{\exp s_{j}}{\sum_{j'=1}^{k} \exp s_{j}'}$$

Interpretation

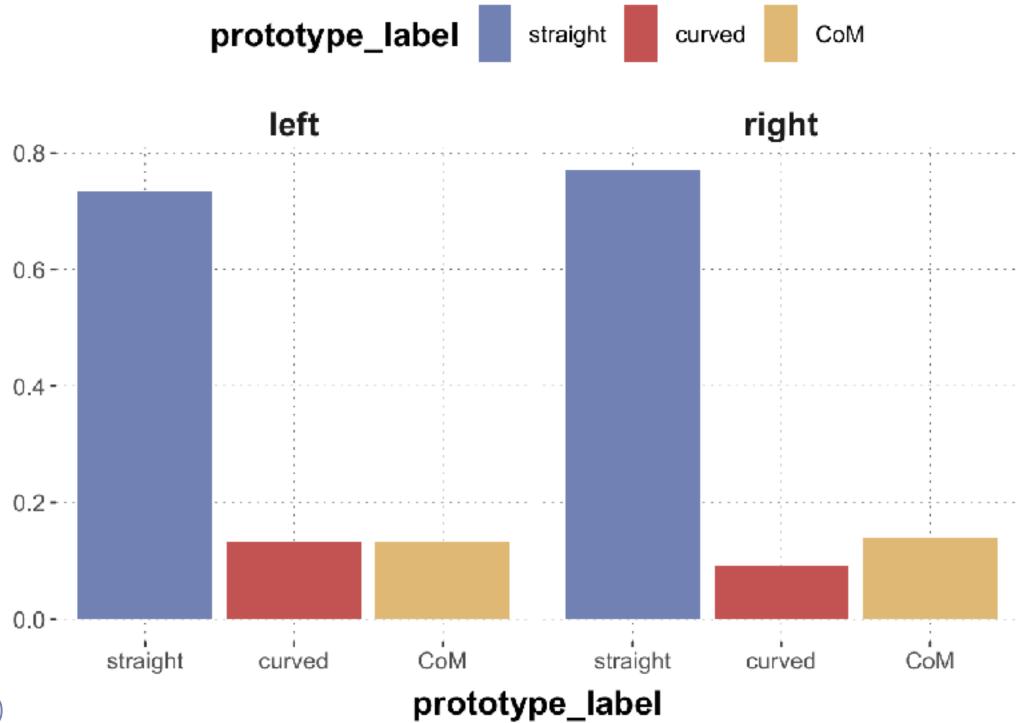
- In linear predictor predictor s_j represents $\log(p_i/p_1)$
- Think of this as ...

k-1 parallel logistic regressions

```
fit_multinom <- brm(
  formula = prototype_label ~ target_position,
  data = data_MT_prepped,
  family = categorical()
)</pre>
```

Multinomial regression

```
fit_multinom <- brm(
  formula = prototype_label ~ target_position,
  data = data_MT_prepped,
  family = categorical()
)</pre>
```



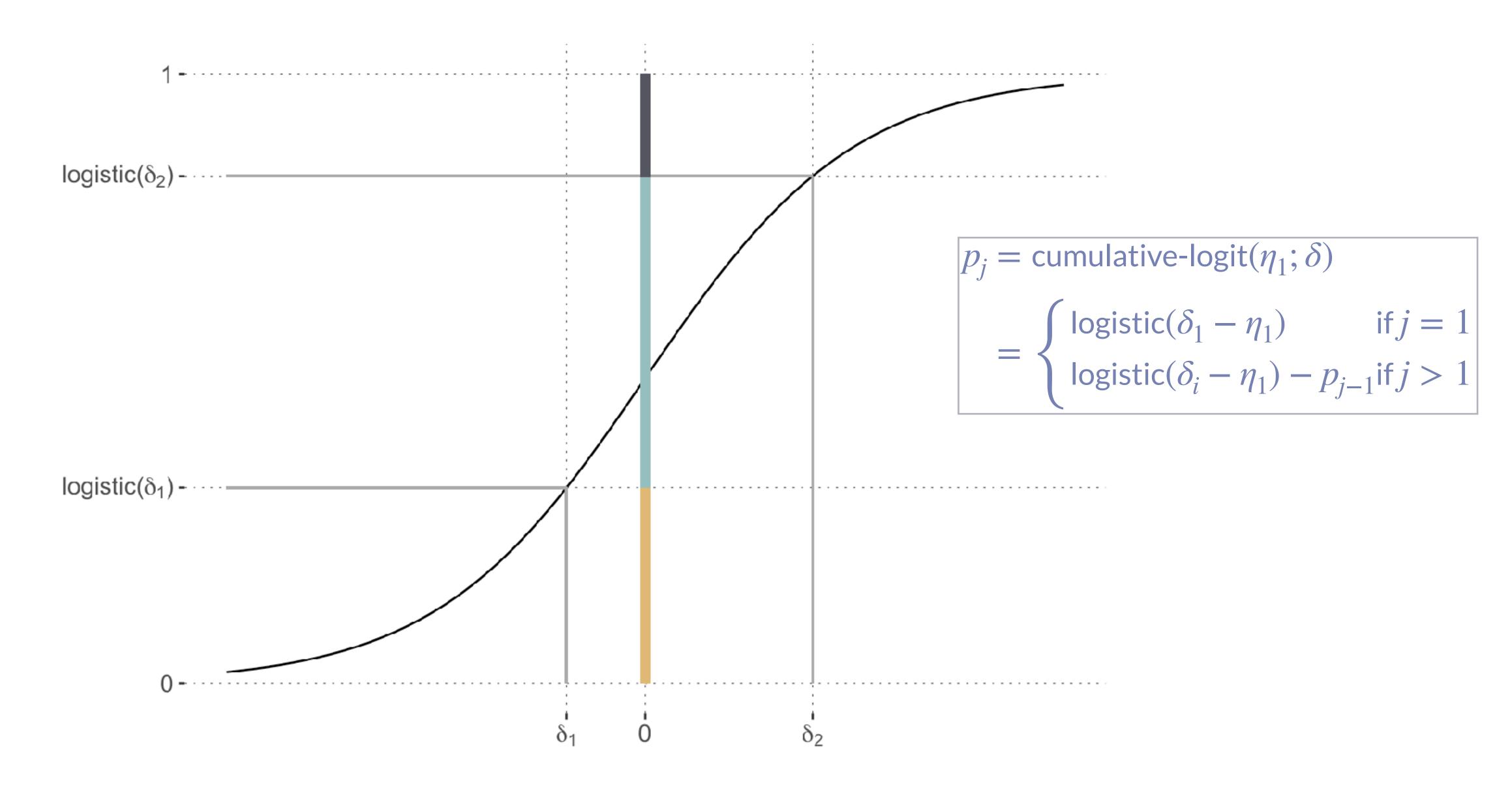
```
Family: categorical
Links: mucurved = logit; muCoM = logit
Formula: prototype_label ~ target_position
  Data: data_MT_prepped (Number of observations: 2052)
  Draws: 4 chains, each with iter = 2000; warmup = 1000; thin = 1;
  total post-warmup draws = 4000
```

Population-Level Effects:

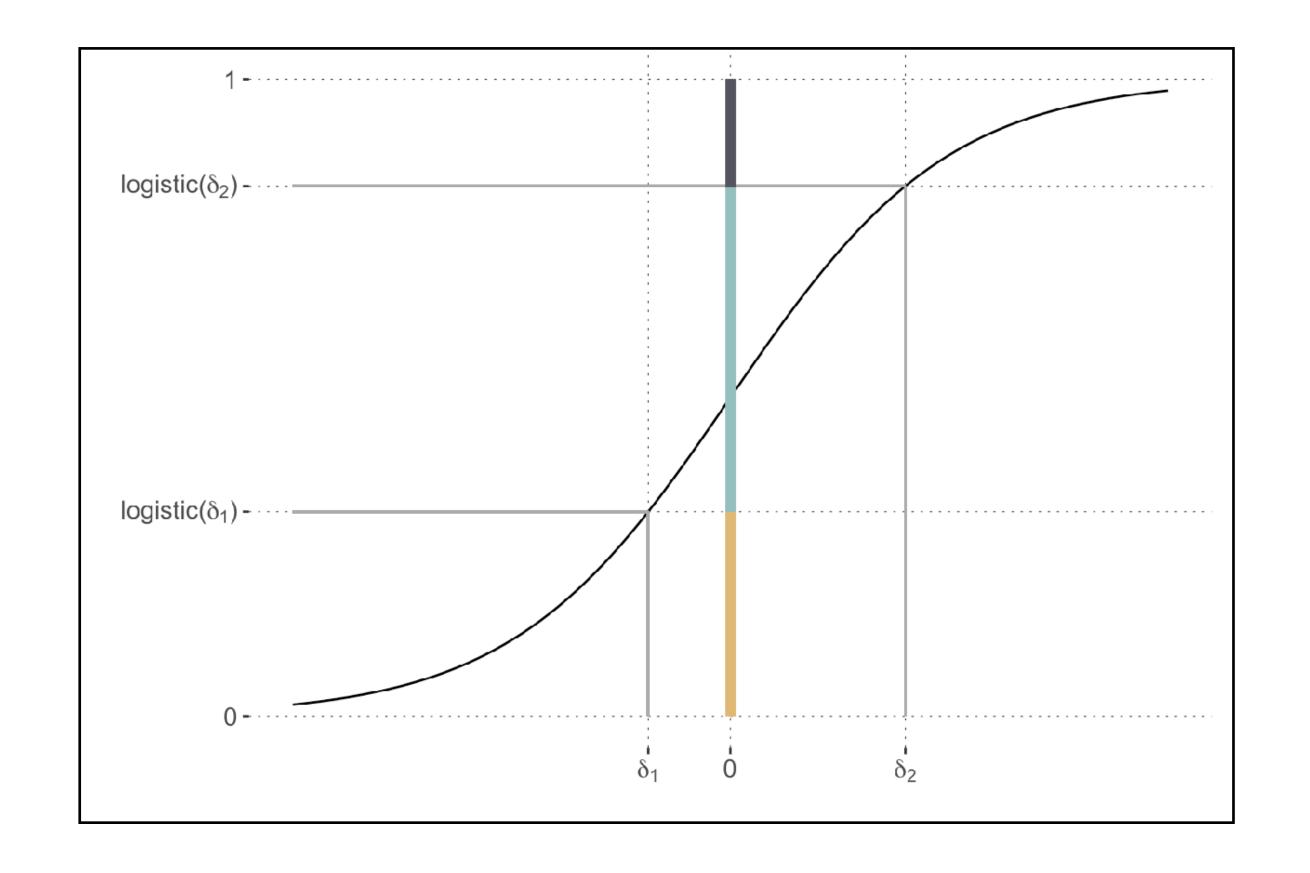
	Estimate	Est.Error	l-95% CI	u-95% CI Rhat
mucurved_Intercept	-1.71	0.09	-1.89	-1.53 1.00
muCoM_Intercept	-1.71	0.09	-1.90	-1.53 1.00
<pre>mucurved_target_positionright</pre>	-0.44	0.14	-0.72	-0.16 1.00
<pre>muCoM_target_positionright</pre>	-0.00	0.13	-0.26	0.25 1.00
	Bulk_ESS	Tail_ESS		
mucurved_Intercept	4595	3205		
muCoM_Intercept	3696	2886		
<pre>mucurved_target_positionright</pre>	3925	3178		
<pre>muCoM_target_positionright</pre>	3910	3066		

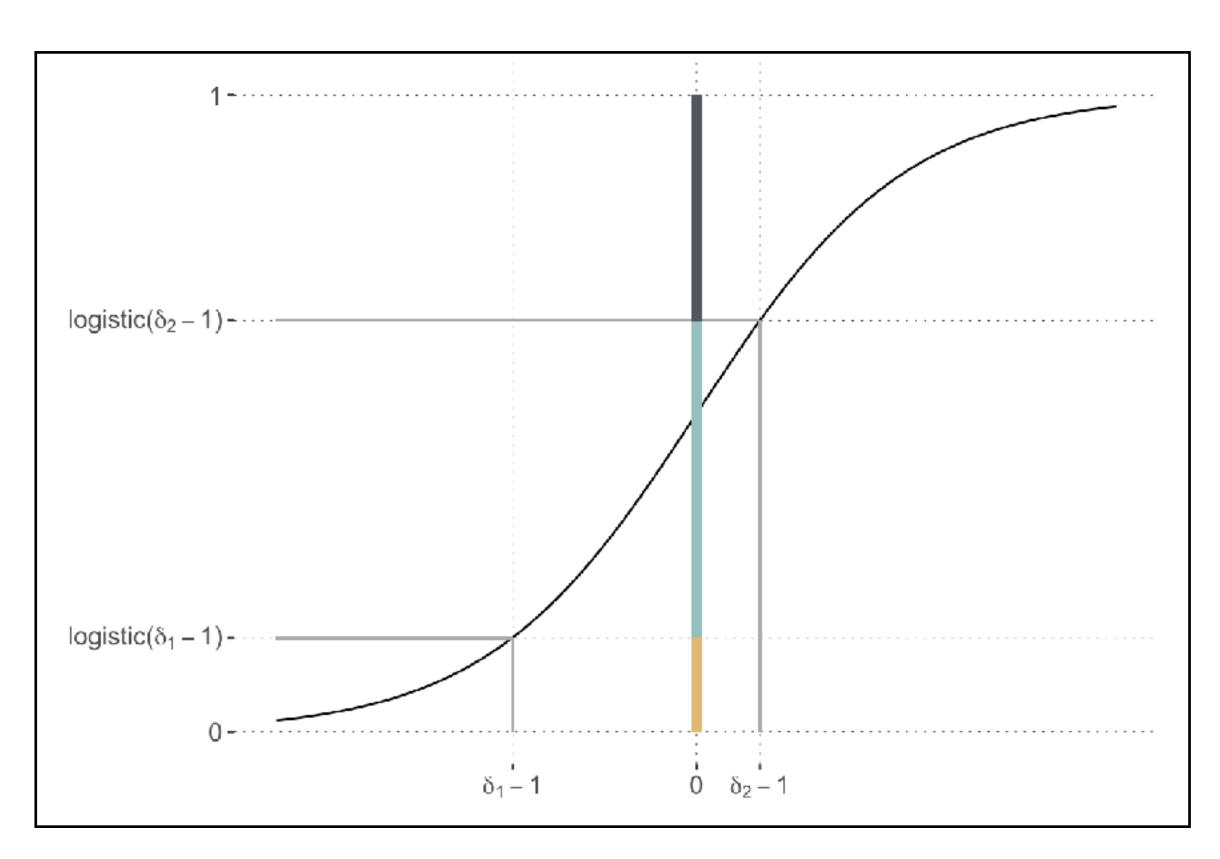
Draws were sampled using sampling(NUTS). For each parameter, Bulk_ESS and Tail_ESS are effective sample size measures, and Rhat is the potential scale reduction factor on split chains (at convergence, Rhat = 1).

cumulative logit



cumulative logit





$$\begin{split} \eta_i &= \mathbf{x}_i \cdot \boldsymbol{\beta} \\ \xi_i &= \text{cumulative-logit}(\eta_i; \boldsymbol{\delta}) \\ y_i &\sim \text{Categorical}(\xi_i) \end{split}$$

[linear predictor]
[predictor of central tendency]
[likelihood]

cumulative logit

```
fit_ordinal <- brm(
  formula = prototype_label ~ MAD,
  data = data_MT_prepped2,
  family = cumulative()
)</pre>
```

Population-Level Effects:

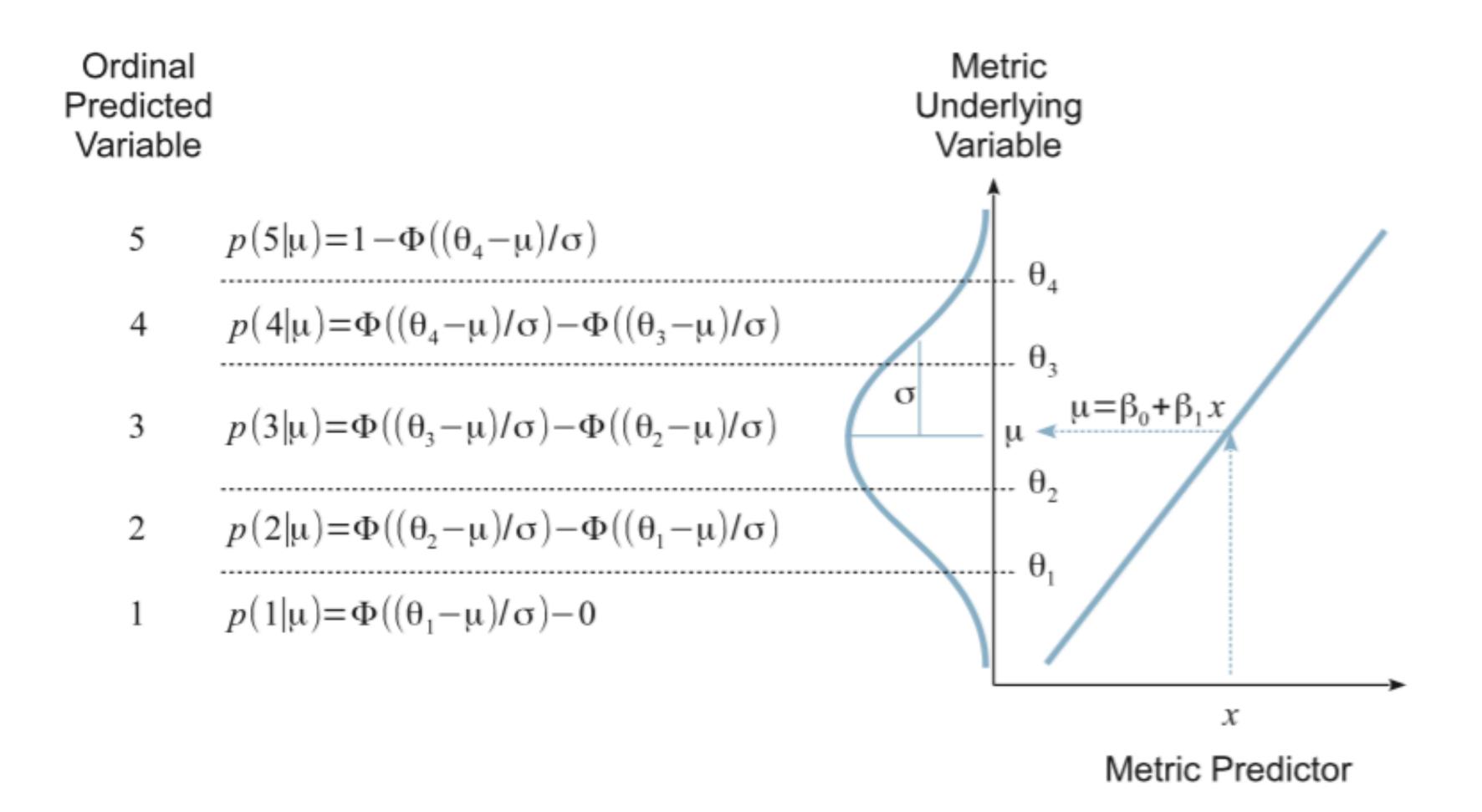
```
Estimate Est.Error l-95% CI u-95% CI Rhat Bulk_ESS Tail_ESS
Intercept[1]
                4.05
                          0.18
                                            4.40 1.00
                                                          2216
                                                                   2247
                                   3.71
Intercept[2]
                9.52
                          0.51
                                   8.56
                                          10.54 1.00
                                                          2003
                                                                  1979
                0.02
                          0.00
                                   0.02
                                            0.03 1.00
                                                          2543
                                                                   2514
MAD
```

Family Specific Parameters:

```
Estimate Est.Error l-95% CI u-95% CI Rhat Bulk_ESS Tail_ESS disc 1.00 0.00 1.00 1.00 NA NA NA
```

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).

cumulative probit



one happy families

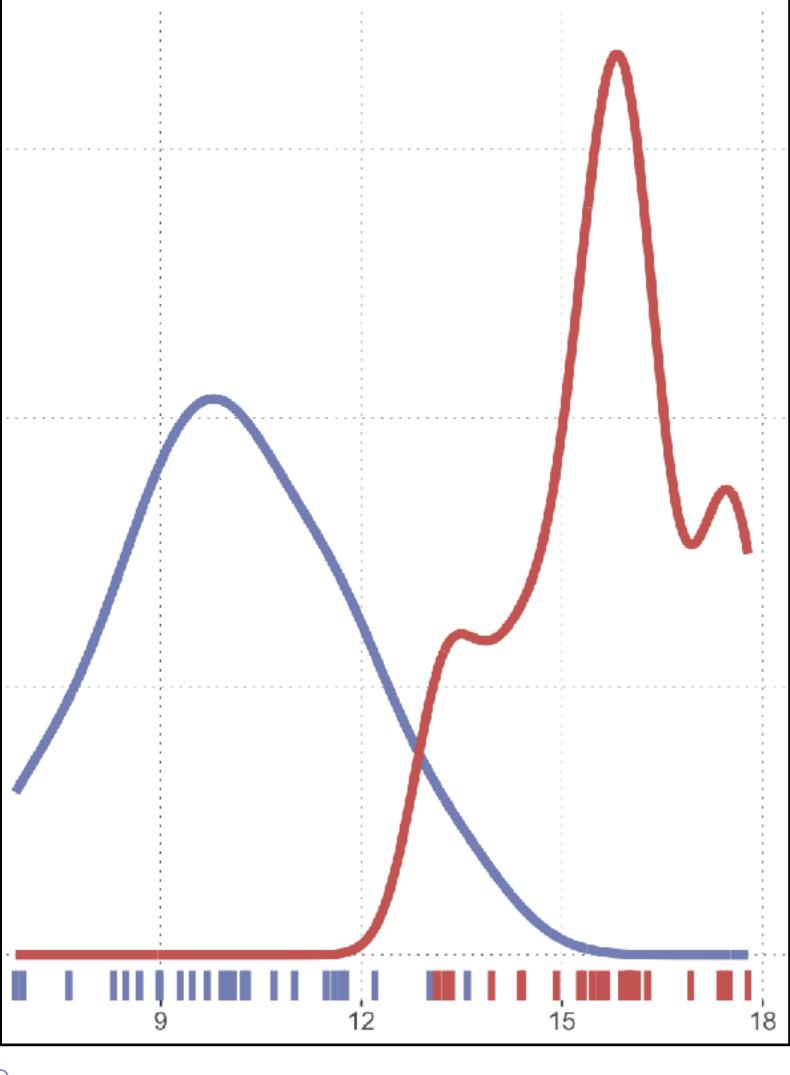
Families in BRMS

- list of available families: <u>link</u>
- explanation of available families: <u>link</u>
- how to write you own: link

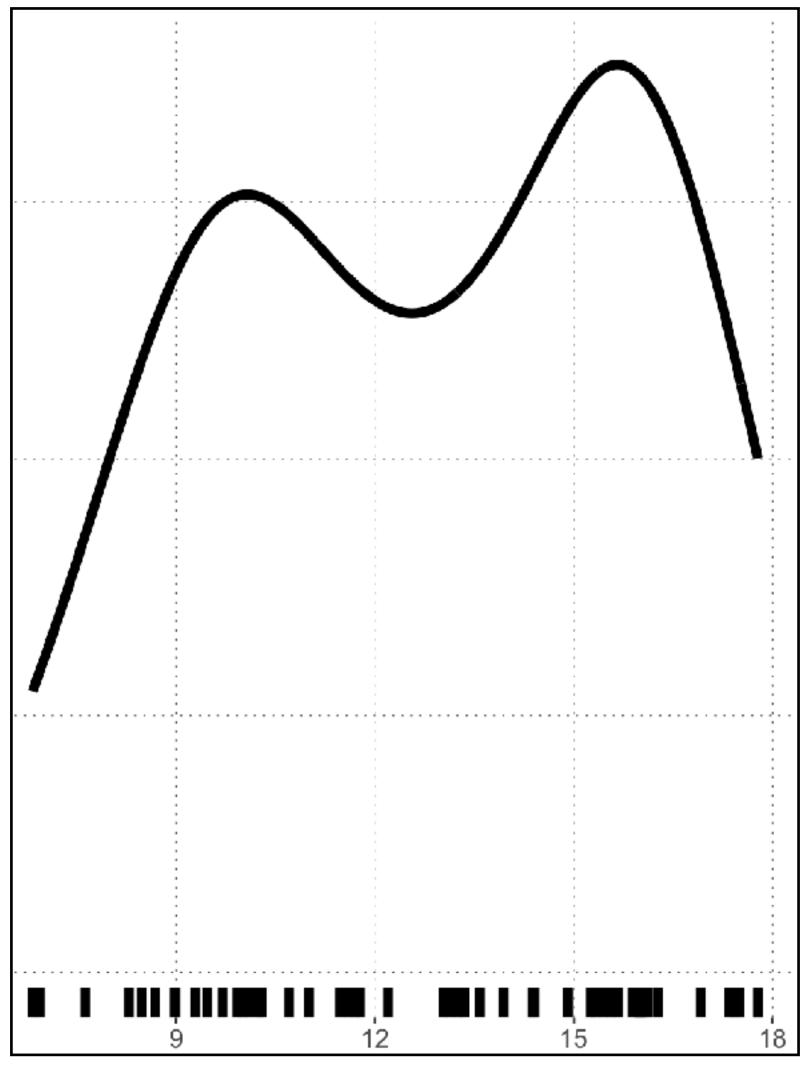
finite mixture models

Multi-modal response distributions

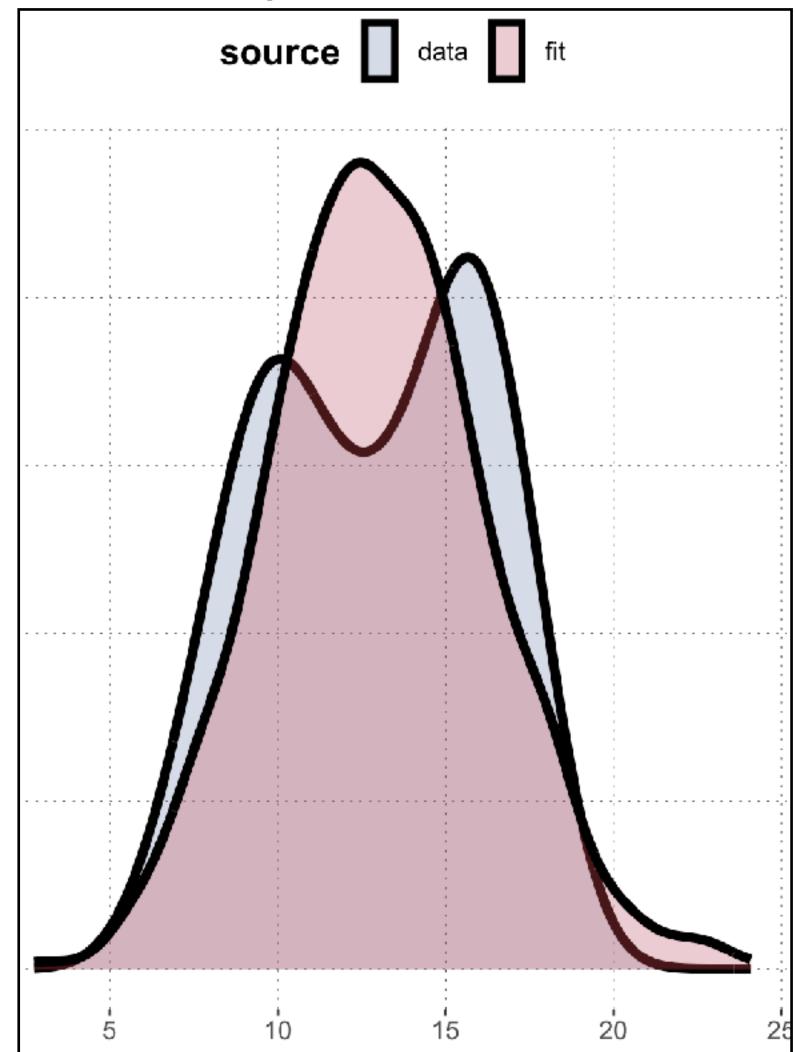
Two-component data



Treated as one for analysis



Posterior predictive ... failure



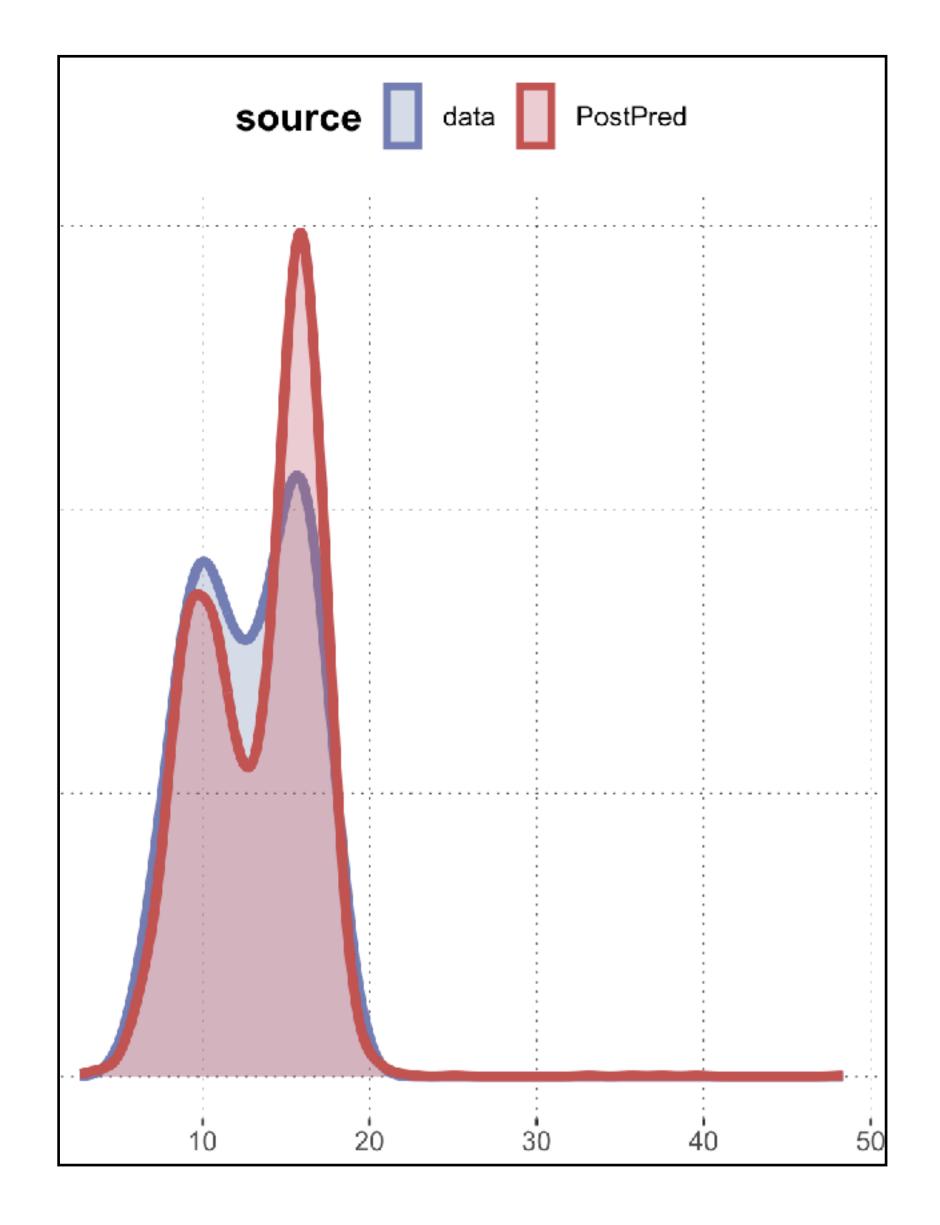
Mixture models

mixing multiple components in the LH function

Let $\langle f_1, ..., f_k \rangle$ be k likelihood functions for data Y. The k-mixture model for Y explains the data as a weighted combination, with mixture weights α (a probability vector). The mixture likelihood function is:

$$f^{\text{MM}}(y_i) = \alpha_{k(i)} f_{k(i)}$$

where k(i) is the mixture component associated with observation i.



Gaussian mixture models

in BRMS

```
brms_fit_2e_GMM <- brm(</pre>
 # intercept only model
  formula = y \sim 1,
  data = data_GMM,
 # declare that the likelihood should be a mixture
 family = mixture(gaussian, gaussian),
 # use weakly informative priors on mu
 prior = c(
    prior(normal(12, 10), Intercept, dpar = mu1),
    prior(normal(12, 10), Intercept, dpar = mu2)
```

special syntax for mixture LH

one intercept for each component

Gaussian mixture models

in BRMS

Population-Level Effects:

mu1_Intercept	10.43	1.25	8.84	12.74	
mu2_Intercept	15.56	0.80	13.28	16.66	

means estimated for each component

Family Specific Parameters:

Estimate Est.Error l-95% CI u-95% CI

		Localide	ESCILITO	C 33 0 CI	u 330 CI	
	sigma1	2.23	0.76	1.16	3.66	
	sigma2	1.59	0.85	0.66	3.31	
Ī						
	theta1	0.54	0.16	0.18	0.88	
	theta2	0.46	0.16	0.12	0.82	

SDs estimated for each component

estimated weights of each component

Zero-inflation models

zeros can be generated by two independent paths

If *f* is a likelihood function for data *y*, the **zero-inflated** likelihood function is:

$$f^{0\inf}(y; \theta, z) = \begin{cases} z + (1 - z) f(y; \theta) & \text{if } y = 0\\ (1 - z) f(y; \theta) & \text{otherwise} \end{cases}$$

both mixture components contribute to likelihood of "zero"

Zero-hurdle models

zeros are generated by one independent path

If *f* is a likelihood function for data *y*, the **zero-hurdle** likelihood function is:

$$f^{0\text{hur}}(y; \theta, z) = \begin{cases} z & \text{if } y = 0 \\ (1 - z) \frac{f(y; \theta)}{1 - f(0; \theta)} & \text{otherwise} \end{cases}$$
only one mixture component contributes to likelihood of "zero"; the other is truncated

Zero/one-inflation models

zeros and can be generated by independent paths

If *f* is a likelihood function for data *y*, the **zero/one-inflated likelihood function** is:

$$f^{0/1\inf}(y; \theta, \alpha, \beta) = \begin{cases} \alpha \gamma & \text{if } y = 1\\ \alpha (1 - \gamma) & \text{if } y = 0\\ (1 - \alpha) f(y; \theta) & \text{otherwise} \end{cases}$$

Hurdle and inflation models in BRMS

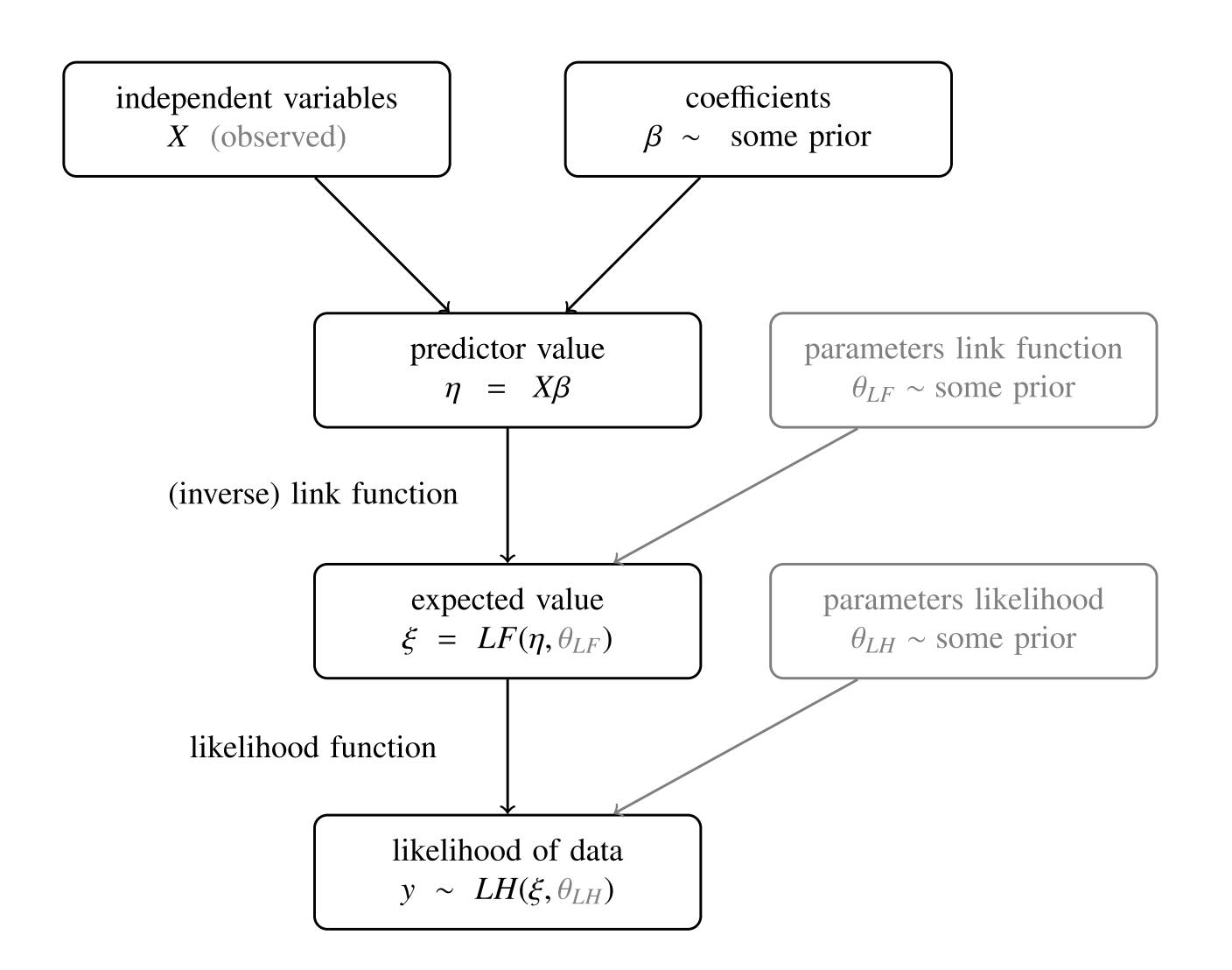
as of March 2025

```
hurdle_poisson(link = "log", link_hu = "logit")
hurdle_negbinomial(link = "log", link_shape = "log", link_hu = "logit")
hurdle_gamma(link = "log", link_shape = "log", link_hu = "logit")
hurdle_lognormal(link = "identity", link_sigma = "log", link_hu = "logit")
hurdle_cumulative(
 link = "logit",
  link_hu = "logit",
  link_disc = "log",
  threshold = "flexible"
zero_inflated_beta(link = "logit", link_phi = "log", link_zi = "logit")
zero_one_inflated_beta(
 link = "logit",
 link_phi = "log",
  link_zoi = "logit",
  link_coi = "logit"
zero_inflated_poisson(link = "log", link_zi = "logit")
zero_inflated_negbinomial(link = "log", link_shape = "log", link_zi = "logit")
zero_inflated_binomial(link = "logit", link_zi = "logit")
zero_inflated_beta_binomial(
  link = "logit",
 link_phi = "log",
  link_zi = "logit"
```

distributional models

Distributional models

linear predictors for parameters of the link- and likelihood functions



Normal GLM:

$$\eta = X\beta$$

$$\theta_{LF}, \theta_{LH} \sim \text{some prior}$$

Distributional GLM:

$$\eta = X\beta$$

$$\theta_{LF} = F(X'\beta')$$

$$\theta_{LF} = F(X''\beta'')$$



recap & preparation

Recap & preparation

▶ recap

- multiple regression & categorical predictors
- generalized linear regression (& beyond)

preparation

- multi-level models
- model comparison
- model criticism