

# Default Priors and Robust Estimation for Generalized Linear Models

(A.K.A., A few things I learned from Luis Pericchi)

---

**Abel Rodríguez**

O'Bayes Meeting  
April, 2022

## Outline

---

Once upon a time ...

Training samples

Prior matching

Heavy tail priors

Concluding remarks

Once upon a time ...

Training samples

Prior matching

Heavy tail priors

Concluding remarks



Once upon a time ...

**Training samples**

Prior matching

Heavy tail priors

Concluding remarks

## “Training” your improper prior

- The idea of splitting your sample into a training set and a test set is at the core of statistics and machine learning.
- Use the training set,  $y^*$ , to turn your improper prior into a proper one, and the rest,  $y^{**}$ , to run the actual test.

$$\pi_k^I(\theta_k | M = k, y^*) = \frac{p_k(y^* | \theta_k, M = k)\pi_k^N(\theta_k | M = k)}{m_k^N(y^*)}$$

$$m_k^I(y^{**} | M = k, y^*) = \int p_k(y^{**} | y^*, \theta_k, M = k)\pi_k^I(\theta_k | M = k, y^*)d\theta_k$$

$$B_{k,k'}^I(y) = \frac{m_k^I(y^{**} | M = k, y^*)}{m_{k'}^I(y^{**} | M = k', y^*)}$$

- If you average over training samples  $\Rightarrow$  *Intrinsic Bayes Factor* (Berger & Pericchi, 1996).
- Size of the training sample?  $\Rightarrow$  **Minimal!**

## EP priors

---

- Introduced by Pérez & Berger (2002).
- EP priors follow a similar rationale, but use an *imaginary* training sample, which is averaged out!

$$\pi_k^E(\boldsymbol{\theta}_k | M_k) = \int \frac{p_k(y^* | \boldsymbol{\theta}_k, M = k) \pi_k^N(\boldsymbol{\theta}_k | M = k)}{m_k^N(y^*)} m^*(y^*) dy^*$$

- Choosing the training sample is replaced with choosing  $m^*$ .
- For nested models, a common option is  $m^*(y^*) = m_0^N(y^*)$  ("simplest" model), which makes the EP prior asymptotically equivalent to the prior implied by the AIBF!

## PEP priors

---

- Introduced by Fouskakis et al. 2015.
- Similar to the EP prior, but it “scales” training likelihood to be (approximately) unit information!

$$\pi_k^{PEP}(\boldsymbol{\theta}_k | M_k) \propto \int \frac{\{p_k(y^* | \boldsymbol{\theta}_k, M = k)\}^{1/\delta} \pi_k^N(\boldsymbol{\theta}_k | M = k)}{\int \{p_k(y^* | \boldsymbol{\theta}_k, M = k)\}^{1/\delta} \pi_k^N(\boldsymbol{\theta}_k | M = k) d\boldsymbol{\theta}_k} m^*(y^* | \delta) p(\delta) dy^* d\delta$$

where  $p(\delta)$  has mean  $n^*$ .

- Computationally tractable when  $p_k(y^* | \boldsymbol{\theta}_k, M = k)$  is Gaussian.



## An example: Linear Models

---

- When  $y \mid \boldsymbol{\theta}_k, \sigma^2, \mathbf{X}_k \sim \mathbf{N}(y \mid \mathbf{X}_k \boldsymbol{\theta}_k, \sigma^2 I)$  and  $\pi_k^N(\boldsymbol{\theta}_k) \propto \mathbf{1}$  and

$$\pi_k^{PEP}(\boldsymbol{\theta}_k \mid M_k) = \int \mathbf{N}\left(\boldsymbol{\theta}_k \mid \left\{\mathbf{X}_k^{*T} \mathbf{X}_k^*\right\}^{-1} \mathbf{X}_k^{*T} \mathbf{y}^*, \delta \sigma^2 \left\{\mathbf{X}_k^{*T} \mathbf{X}_k^*\right\}^{-1}\right) m^*(\mathbf{y}^* \mid \delta) p(\delta) d\delta d\mathbf{y}^*$$

- For  $n^* = n$  and  $\mathbf{X}_k^* = \mathbf{X}_k$ , compare that with the corresponding g-prior:

$$\pi_k^g(\boldsymbol{\theta}_k \mid M_k) = \int \mathbf{N}\left(\boldsymbol{\theta}_k \mid \mathbf{0}, \delta \sigma^2 \left\{\mathbf{X}_k^{*T} \mathbf{X}_k^*\right\}^{-1}\right) \tilde{p}(\delta) d\delta$$

(Recall that  $\tilde{p}$  is centered around  $n$ .)

## Generalizing PEPs to GLMs

---

- In the case of Gaussian linear models, the PEP is relatively easy to derive because rescaling by  $1/\delta$  leaves the likelihood in the normal family.
- The same is not true for other members of the exponential family (e.g., logistic or loglinear regression).
- Fouskakis et al. (2018) propose a generalization, but it has various theoretical, computational and empirical drawbacks.
- We propose a different generalization: the Laplace PEPs!
  - ▶ Porwal, A., & Rodriguez, A. (2021). Laplace Power-expected-posterior priors for generalized linear models with applications to logistic regression. arXiv preprint arXiv:2112.02524.

## Laplace Power-expected-posterior Prior (LPEP)

---

- To construct the Laplace PEP, replace  $p_k(y^* | \theta_k, M = k)$  with its Laplace approximation **before** raising to the  $1/\delta$  power!

$$\pi_k^{PEP}(\theta_k | M_k) \propto \int \frac{\mathbf{N}(\theta_k | \hat{\theta}_k(y^*), \delta H_k^{-1}(\hat{\theta}_k(y^*))) \pi_k^N(\theta_k | M = k)}{\int \mathbf{N}(\theta_k | \hat{\theta}_k(y^*), \delta H_k^{-1}(\hat{\theta}_k(y^*))) \pi_k^N(\theta_k | M = k) d\theta_k} m^*(y^*) p(\delta) dy^* d\delta$$

- Laplace approximation should be particularly accurate when  $n^* = n$  but, conceptually, the procedure is reasonable for other choices of  $n^*$ .
- An implicit constraint on  $y^*$  is that the Laplace approximation needs to be well defined (e.g., the MLE needs to exist for every model under consideration).

## Example: Logistic regression

---

- Likelihood is

$$p_k(y \mid \theta_k, M = k) = \prod_{i=1}^n \frac{\exp \{ y_i x_{i,k}^T \theta_k \}}{1 + \exp \{ x_{i,k}^T \theta_k \}}$$

- Pick  $m^*(y^*) = m_0^N(y^*) 1(y^* \in \Omega_k(X))$  where

$$\Omega_k(X) = \{y : \hat{\theta}_k(y, X) \text{ is finite for all } k\}$$

and

$$m_0^N(y^*) = \frac{\Gamma(y^* + 1/2) \Gamma(n - y^* + 1/2)}{\Gamma(n + 1) (\Gamma(1/2))^2}$$

- Various possible choices for  $p(\delta)$  (fixed, hyper-g/n, robust prior).

## Example: Logistic regression

---

- For logistic regression (and many other GLMs!) it is enough to show that  $\hat{\theta}_k(y, X)$  is finite when  $k$  corresponds to the full model!
  - ▶ We provide easy-to-verify sufficient conditions in the paper.
- Checking this condition for logistic regression is relatively straightforward using the algorithm of Kosmidis and Schumacher (2020).
- In this case, the prior is proper for every model  $k$ .

## Properties of the LPEP for GLMs

---

- For linear models, this is just your standard PEP!
- Under standard regularity conditions Bayes factors / posterior model probabilities are consistent.
  - ▶ True even if  $p$  grows with  $n$  at a reasonably slow rate.
- Well-defined intrinsic prior.
- Unlike Li & Clyde, (2018), it can be used even if the original data is separable, or in hierarchical settings.
  - ▶ Good theoretical properties.
- The fact that they correspond to mixtures of normals facilitates computation using MCMC.
  - ▶ For a number of GLMs, no need for reversible-jump schemes like Fouskakis et al. (2018).
  - ▶ A second Laplace approximation can be used to speed up computation as in Li & Clyde, (2018).

## Simulation study: Design

- $n = 500$ ;  $p = 100$ ; 100 bootstrapped datasets
- Columns of  $X$  drawn from standard normal distribution with pairwise correlation  $\text{cor}(X_i, X_j) = r^{|i-j|}$  for  $1 \leq i < j \leq p$
- Scenarios:  $r = 0$  (independent design) and  $r = 0.75$  (correlated design)
- $p_{M_T}$  denote the number of variables in the true model
- $\mathbf{b} = (2, -1, -1, 0.5, -0.5)^T$  and  $\beta_{M_T, 21:100} = 0$

$p_{M_T}$	$\beta_{M_T,0}$	$\beta_{M_T,1:5}$	$\beta_{M_T,6:10}$	$\beta_{M_T,11:15}$	$\beta_{M_T,16:20}$
0	-0.5	0	0	0	0
5	-0.5	$\mathbf{b}$	0	0	0
10	-0.5	$\mathbf{b}$	0	$\mathbf{b}$	0
20	-0.5	$\mathbf{b}$	$0.5\mathbf{b}$	$\mathbf{b}$	$0.5\mathbf{b}$

- Comparison with: Mixture of g-priors (Li & Clyde, 2018), LASSO, SCAD and MCP.

## Simulation study: Results - MAP model properties

$p$		100							
		Beta-Binomial(1,1)							
$p(M)$		0		5		10		20	
$pM_T$		0		5		10		20	
$r$		0	0.75	0	0.75	0	0.75	0	0.75
$\delta = n$	LPEP	99	<b>100*</b>	45	4	<b>18*</b>	0	0	0
	LCE	<b>100*</b>	<b>100*</b>	45	<b>5</b>	8	0	0	0
	LCL	<b>100*</b>	<b>100*</b>	<b>46</b>	4	11	0	0	0
$\delta \sim \text{robust}$	LPEP	99	<b>100*</b>	<b>53*</b>	<b>6*</b>	<b>15</b>	0	0	0
	LCE	99	<b>100*</b>	45	<b>6*</b>	0	0	0	0
	LCL	<b>100*</b>	<b>100*</b>	46	<b>6*</b>	2	0	0	0
$\delta \sim \text{hyper } g/n$	LPEP	<b>98</b>	<b>100*</b>	<b>50</b>	<b>5</b>	<b>17</b>	0	0	0
	LCE	97	99	25	4	0	0	0	0
	LCL	65	78	3	0	0	0	0	0
	LASSO	59	65	0	0	0	0	0	0
	SCAD	57	59	0	0	0	0	0	0
	MCP	<b>73</b>	<b>66</b>	<b>8</b>	0	<b>3</b>	0	0	0

**Table:** Number of times (**over 100 replications**) that the **MAP** model coincides with the true model in the logistic regression ; **BOLD** represent group maximum; \* represent overall maximum.



# Simulation study: Results - F1 score

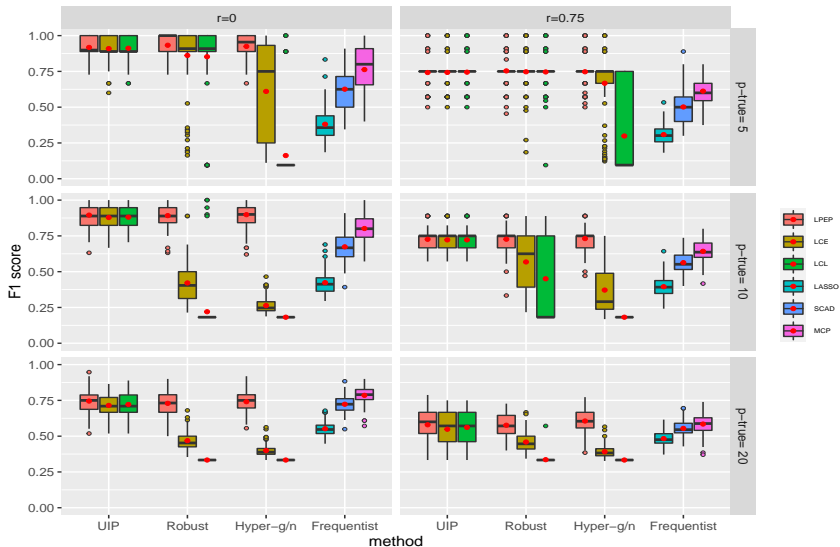


Figure: F1 score across 100 simulated datasets; Red dots represent the average

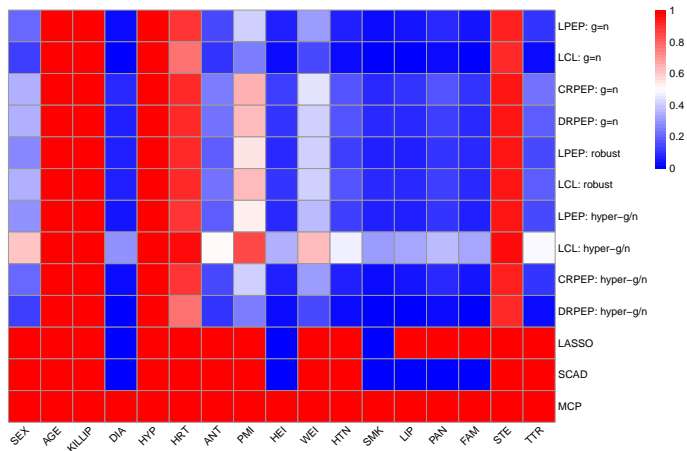
## Simulation study: Results - Mean squared error

$p$		100							
		Beta-Binomial(1,1)							
$p(M)$		0		5		10		20	
$pM_T$		0		5		10		20	
$r$		0	0.75	0	0.75	0	0.75	0	0.75
$\delta = n$	LPEP	0.11	<b>0.10*</b>	2.91	<b>7.67</b>	7.09	<b>17.67</b>	<b>14.70</b>	<b>33.90</b>
	LCE	0.11	<b>0.10*</b>	3.06	7.78	7.64	18.44	16.11	36.47
	LCL	<b>0.10*</b>	<b>0.10*</b>	<b>2.87</b>	7.68	<b>6.78</b>	18.17	16.22	36.43
$\delta \sim \text{robust}$	LPEP	0.12	<b>0.10*</b>	<b>2.62*</b>	<b>6.87*</b>	<b>6.04*</b>	<b>14.07*</b>	<b>13.38</b>	<b>24.03*</b>
	LCE	0.12	0.11	4.83	7.80	47.30	23.30	96.14	52.80
	LCL	<b>0.10*</b>	<b>0.10*</b>	8.86	8.44	214.63	60.56	275.93	115.58
$\delta \sim \text{hyper g/n}$	LPEP	<b>0.16</b>	0.14	<b>2.70</b>	<b>6.89</b>	<b>6.12</b>	<b>14.76</b>	<b>13.03*</b>	<b>24.86</b>
	LCE	0.23	<b>0.13</b>	6.71	8.90	38.54	26.25	51.48	44.29
	LCL	0.29	0.31	34.28	22.93	104.10	72.95	130.80	94.98
	LASSO	0.25	0.20	7.08	11.91	17.15	25.04	29.44	36.69
	SCAD	<b>0.21</b>	<b>0.16</b>	3.07	9.02	6.62	<b>18.80</b>	14.88	<b>33.00</b>
	MCP	0.22	0.18	<b>2.82</b>	<b>8.92</b>	<b>6.35</b>	19.38	<b>15.13</b>	33.52

**Table:** 1000 times the AMSE for estimated coefficients over 100 replications; **BOLD** represent group minimum; \* represent overall minimum.

## Gusto-I study: survival to treatments for occluded coronary arteries

- Model the binary endpoint of 30-day survival for a subgroup of  $n = 2188$  patients using 17 clinical covariates



**Figure:** Marginal posterior inclusion probabilities (PIPs) for GUSTO-I dataset (Bayesian procedures) and variables included in the model (penalized likelihood methods).

## GUSTO-I study: Out-of-sample Predictive Performance

- We performed a 10-fold cross-validation study
- AUC and Calibration slope (CS) allow us to evaluate the methods in terms of discrimination and calibration;  $\uparrow$  score is better
- LS & BRIER measure the predictive accuracy of methods;  $\downarrow$  score is better

		AUC	CS	LS	BRIER
$\delta = n$	LPEP	<b>0.8324*</b>	0.9971	<b>0.1824</b>	<b>0.0496</b>
	LCL	0.8300	0.9931	0.1831	0.0497
	CRPEP	0.7789	1.0578	0.1965	0.0521
	DRPEP	0.7790	1.0569	0.1963	0.0521
$\delta \sim$ robust	LPEP	<b>0.8322</b>	<b>1.0129</b>	<b>0.1822</b>	<b>0.0495</b>
	LCL	0.8316	0.9804	<b>0.1822</b>	<b>0.0495</b>
$\delta \sim$ hyper g/n	LPEP	<b>0.8319</b>	<b>1.0074*</b>	0.1823	0.0495
	LCL	0.8311	1.0109	<b>0.1818</b>	<b>0.0493</b>
	CRPEP	0.7956	1.1677	0.1951	0.0522
	DRPEP	0.7800	1.0571	0.1961	0.0520
	LASSO	<b>0.8305</b>	<b>1.0369</b>	<b>0.1816*</b>	<b>0.0492*</b>
	SCAD	0.8243	0.9135	0.1838	0.0496
	MCP	0.8250	0.9196	0.1838	0.0496

**Table:** Average prediction accuracy measures in a 10-fold cross validation study for GUSTO-I dataset 20 / 44

Once upon a time ...

Training samples

**Prior matching**

Heavy tail priors

Concluding remarks

## Factor models

---

- Consider multivariate responses  $y_i = (y_{i,1}, \dots, y_{i,J})^T$  where  $y_{i,j} \in \mathbb{R}$  and  $i = 1, \dots, I$ . A factor model takes the form

$$y_{i,j} = \mu_j + \alpha_j^T \beta_i + \epsilon_{i,j} \quad \epsilon_{i,j} \sim \text{N}(0, \sigma_j^2)$$

$$\alpha_j^T = (\alpha_{j,1}, \dots, \alpha_{j,d}), \beta_i^T = (\beta_{i,1}, \dots, \beta_{i,d}), \text{ and } d \ll J.$$

- Used for dimensionality reduction, covariance estimation, prediction.
- The same bilinear structure can be built into Generalized Linear Models. For example, for binary data  $y_{i,j} \in \{0, 1\}$ ,

$$y_{i,j} \sim \text{Ber}(\theta_{i,j}) \quad \theta_{i,j} = G_j \left( \mu_j + \alpha_j^T \beta_i \right)$$

where  $G$  is a link function (probit, logit, etc).

- Can be naturally extended to network/relational data.

## Factor models: challenges

---

- Common practical challenges related to model selection:
  - ▶ Selecting the dimension  $d$  of the latent space.
  - ▶ Selecting between a parametric and a non-parametric specification for the distribution of the latent traits.
- The parameters of the model are not identifiable without incorporating some constraints.
  - ▶ This can make interpretation and prior elicitation hard.
- Priors need to be chosen very carefully if comparisons are going to be meaningful.

## Factor models: selecting $d$

---

- Consider a slight generalization of the factor model where

$$y_{i,j} = \mu_j + \boldsymbol{\alpha}_j^T \boldsymbol{\Lambda} \boldsymbol{\beta}_i + \epsilon_{i,j}$$

where  $\boldsymbol{\Lambda} = \text{diag}\{\lambda_1, \dots, \lambda_d\}$  and  $\lambda_k \in \{0, 1\}$

- The introduction of the  $\lambda_k$ s would in principle enable inference of the dimension of the latent space.
- Note that

$$\text{Var}(y_{i,j} \mid \mu_j, \boldsymbol{\alpha}_j, \boldsymbol{\Lambda}) = \text{Var}(\epsilon_{i,j}) + \sum_{k=1}^d \lambda_k \alpha_{j,k} \text{Var}(\beta_{i,k})$$

- If i.i.d. priors are used for the  $\beta_{i,k}$ s (which is common), then

$$\lim_{d \rightarrow \infty} \text{Var}(y_{i,j} \mid \mu_j, \boldsymbol{\alpha}_j, \boldsymbol{\Lambda}) = \infty$$



## Factor models: selecting $d$

---

- There are a couple of possible solutions:
  - ▶ Allow the variance of  $\beta_{i,k}$  to decrease with  $k$  fast enough, for example  $\text{Var}(\beta_{i,k}) = \mathcal{O}(k^{-2})$ .
  - ▶ Allow  $Pr(\lambda_k = 1)$  to decrease fast enough with  $k$ .
- This setting extends to factor models embedded in GLMs.
- We have used these approaches in a few papers:
  - ▶ Guha, S. & Rodriguez, A. (2021). Bayesian regression with undirected network predictors with an application to brain connectome data. *Journal of the American Statistical Association*, 116(534), 581-593.
  - ▶ Sosa, J. & Rodríguez, A. (2021). A latent space model for cognitive social structures data. *Social Networks*, 65, 85-97.
  - ▶ Guhaniyogi, R. & Rodriguez, A. (2020). Joint modeling of longitudinal relational data and exogenous variables. *Bayesian Analysis*, 15(2), 477-503.

**Underlying principle:** when eliciting priors on non-identifiable parameters for various models, the implied priors on key identifiable quantities should be similar across models.

## Factor models: parametric vs. non-parametric specifications

- Consider the 1D factor model:

$$y_{i,j} \sim \text{Ber}(G(\mu_j + \alpha_j \beta_i))$$

- Motivating application: *item response models*

$i$  = Test subject

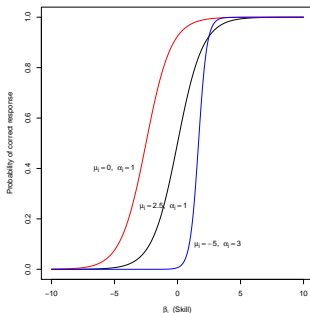
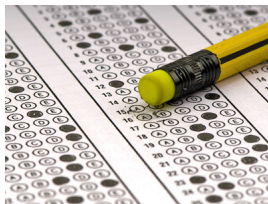
$j$  = Question

$\mu_j$  = Difficulty

$\alpha_j$  = Discrimination

$\beta_i$  = Skill

- Rasch model is a special case.



## Factor models: parametric vs. non-parametric specifications

---

- Two possible specifications for the random effect:
  - ▶ Standard parametric model:  $\beta_i \sim N(0, 1)$
  - ▶ Non-parametric specification (Dirichlet process mixture of normals):

$$\beta_i | G \sim \int N(\cdot | \eta, \tau^2) G(d\eta, d\tau^2), \quad G \sim DP(M, G_0)$$

- How do you fairly compare these two models?
  - ▶ Paganin, S., Paciorek, C. J., Wehrhahn, C., Rodriguez, A., Rabe-Hesketh, S., & de Valpine, P. (2022+). Computational methods for Bayesian semiparametric Item Response Theory models. arXiv preprint arXiv:2101.11583.
  - ▶ Try to match the prior distribution of  $\theta_i = G(\mu_j + \alpha_j \beta_i)$  across both models!

## Binary factor models in general topological spaces

---

- The models we discussed previously project the data on low-dimensional Euclidean spaces.
- In some applications (e.g., in political sciences) other geometries might be more appropriate!

## Spatial voting models

Rational choice theory derivation:

$\psi_j$  = "Yeah" position  $\in \mathbb{R}^d$

$\zeta_j$  = "Nay" position  $\in \mathbb{R}^d$

$\beta_i$  = Ideal point  $\in \mathbb{R}^d$

$$U_{i,j}(\text{Yeah}) = -\|\beta_i - \psi_j\|^2 + \epsilon_{i,j}$$

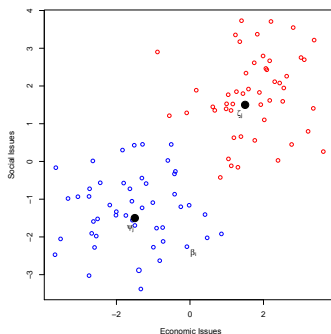
$$U_{i,j}(\text{Nay}) = -\|\beta_i - \zeta_j\|^2 + \nu_{i,j}$$

where  $\nu_{i,j} - \epsilon_{i,j} \sim G_j$ , and  $y_{i,j} = 1 \Leftrightarrow U_{i,j}(\text{Yeah}) > U_{i,j}(\text{Nay})$ ,

$$\mu_j = \zeta_j^T \zeta_j - \psi_j^T \psi_j$$

$$\alpha_j = 2(\psi_j - \zeta_j)$$

Policy Space representation



## Binary factor models in general topological spaces

---

- Consider letting  $\psi_j, \zeta_j, \beta_i \in \mathcal{D}$ , where  $\mathcal{D}$  is a connected Riemannian manifold and define

$$U_{i,j}(\text{Yes}) = - \{d(\beta_i, \psi_j)\}^2 + \epsilon_{i,j},$$

$$U_{i,j}(\text{No}) = - \{d(\beta_i, \zeta_j)\}^2 + \nu_{i,j},$$

where  $d(\beta_i, \psi_j)$  is the *geodesic distance* between  $\beta_i$  and  $\psi_j$  and  $\nu_{i,j} - \epsilon_{i,j} \sim G_{\kappa_j}$ .

- As before,  $y_{i,j} = 1$  iff  $U_{i,j}(\text{Yes}) > U_{i,j}(\text{No})$ , so

$$P(y_{i,j} = 1 \mid \beta_i, \zeta_j, \kappa_j) = G_{\kappa_j} \left( \{d(\beta_i, \zeta_j)\}^2 - \{d(\beta_i, \psi_j)\}^2 \right).$$

## Spherical factor models

---

- In the  $\mathcal{S}^{K+1}$ , the geodesic distance is given by  $\rho_{K+1}(\boldsymbol{\psi}, \boldsymbol{\beta}) = \arccos\left(\mathbf{x}_{\boldsymbol{\psi}}^T \mathbf{x}_{\boldsymbol{\beta}}\right)$ , with, for example,

$$x_{\boldsymbol{\psi},1} = \cos \psi_1 \cos \psi_2 \cos \psi_3 \cdots \cos \psi_{K-1},$$

$$x_{\boldsymbol{\psi},2} = \sin \psi_1 \cos \psi_2 \cos \psi_3 \cdots \cos \psi_{K-1},$$

$$x_{\boldsymbol{\psi},3} = \sin \psi_2 \cos \psi_3 \cdots \cos \psi_{K-1}$$

$$\vdots$$

$$x_{\boldsymbol{\psi},K} = \sin \psi_{K-1} \cos \psi_K,$$

$$x_{\boldsymbol{\psi},K+1} = \sin \psi_K.$$

- Yu, X., & Rodriguez, A. (2022). A Bayesian Approach to Spherical Factor Analysis for Binary Data. arXiv preprint arXiv:2008.05109.

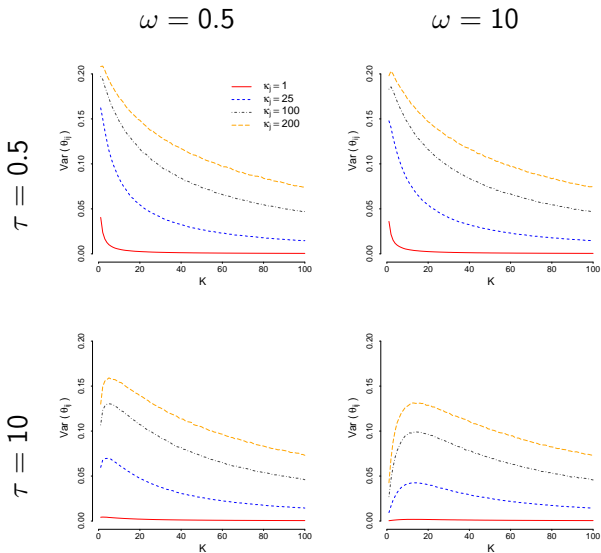


## Priors for spherical factor models

---

Standard von-Mises Fisher distributions on the sphere for  $\{\psi_j\}_{j=1}^J$ ,  
 $\{\zeta_j\}_{j=1}^J$  and  $\{\beta_i\}_{i=1}^I$  will not work!

# Variance of induced prior on $\theta_{i,j}$ - Von Misses-Fisher priors



## Spherical models

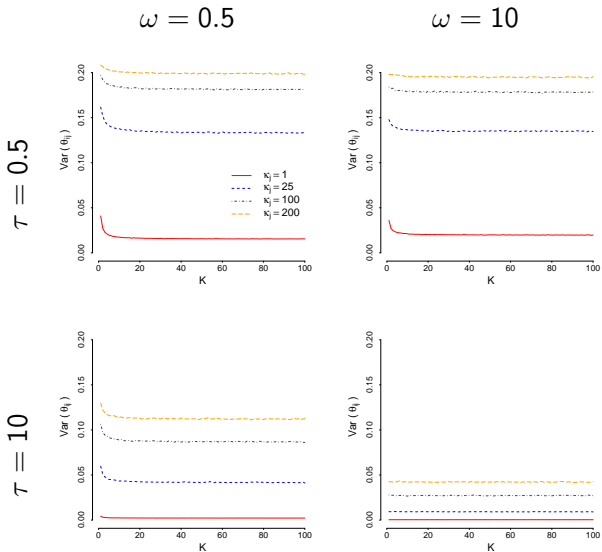
---

- We need a new class of priors on the sphere that allows for marginal variances of the angles to decrease with as new dimensions are added

$$p(\phi | \omega) = \left(\frac{1}{2\pi}\right)^K 2^{K-1} \frac{1}{I_0(\omega_1)} \exp\{\omega_1 \cos \phi_1\} \prod_{k=2}^K \frac{1}{I_0(\omega_k)} \exp\{\omega_k \cos 2\phi_k\}$$

- Unlike the Euclidean case, we need the variance to decrease for both the ideal points and the Yes/No positions!

# Variance of induced prior on $\theta_{i,j}$ - Von Misses-Fisher priors



Once upon a time ...

Training samples

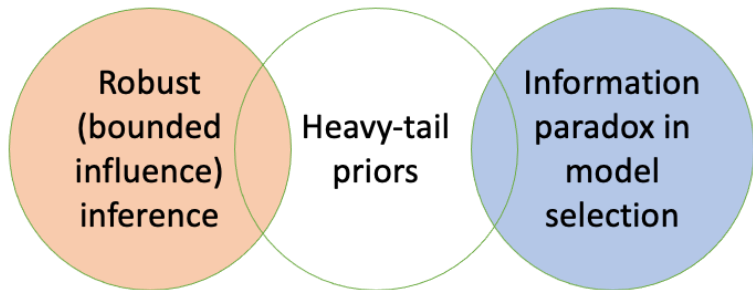
Prior matching

**Heavy tail priors**

Concluding remarks

## The role of priors with heavy tails

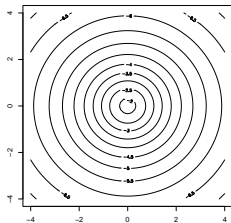
---



## Two philosophies ...

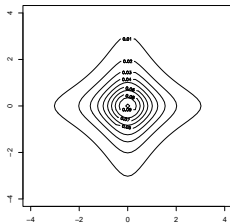
### $g$ -priors and its kin

- Accounts for the “right” correlation among coefficients.
- “Non-directional”: Same tail behavior in every direction



### Horseshoe and its kin

- Coefficients are independent a priori.
- “Directional”: tails along axis are heavier than tails in other directions



## Getting the best of both worlds:

---

- “Directional”  $g$ -priors:

$$\boldsymbol{\theta}_\gamma \mid \gamma, \Lambda_\gamma, \sigma^2 \sim \text{N} \left( 0, \sigma^2 \Lambda_\gamma^{1/2} \left\{ \mathbf{X}_\gamma^T \mathbf{X}_\gamma \right\}^{-1} \Lambda_\gamma^{1/2} \right)$$

with  $\Lambda_\gamma = \text{diag} \{ \lambda_{\gamma,1}, \dots, \lambda_{\gamma,p_\gamma} \}$  and  $\lambda_{\gamma,j} \sim H$ .

- “Correlated” continuous shrinkage priors:

$$\boldsymbol{\theta} \mid \Lambda, \sigma^2 \sim \text{N} \left( 0, \sigma^2 \left\{ \mathbf{X}^T \Lambda^{-1} \mathbf{X} \right\}^{-1} \right)$$

with  $\Lambda = \text{diag} \{ \lambda_1, \dots, \lambda_p \}$  and  $\lambda_j \sim H$ .



## Factor models

---

- A lot of the literature on continuous shrinkage priors has focused on making the Horseshoe a bit more flexible by making the distribution  $H$  more flexible by adding a couple of extra parameters.
- You could make the specification more flexible by setting a non-parametric prior on  $H$  (e.g., a Pòlya Tree centered on the half Cauchy distribution).
- Still somewhat speculative, this is work in progress!
  - ▶ Calibration?
  - ▶ How much can you really learn when you specify a non-parametric model further down in the hierarchy?
  - ▶ Efficient computation.

Once upon a time ...

Training samples

Prior matching

Heavy tail priors

Concluding remarks

## Concluding remarks

---

- A lot of the things that I learned from Pericchi 20 years ago still influence both my research and my teaching.
- I cannot believe it has been 20 years ...
- The school that he created in Venezuela starting in the mid 80s and early 90s is still going strong, if in exile ...

Thank you!