

# The Role of Skewed Distributions in Bayesian Inference

(conjugacy, scalable approximations and asymptotics)

O'Bayes 2022: Objective Bayes Methodology Conference

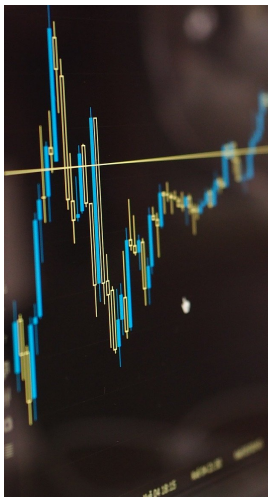
07.09.2022

Daniele Durante

[daniele.durante@unibocconi.it](mailto:daniele.durante@unibocconi.it)

## Regression . . .

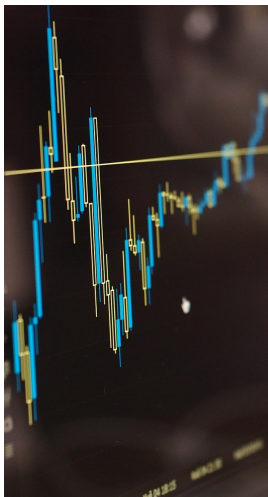
*“Statisticians are engaged in an exhausting but exhilarating struggle with the biggest challenge: **how to translate information into knowledge**” [S. Senn]*



SOURCE: <https://pixabay.com/it/>

# Regression . . .

*“Statisticians are engaged in an exhausting but exhilarating struggle with the biggest challenge: **how to translate information into knowledge**” [S. Senn]*

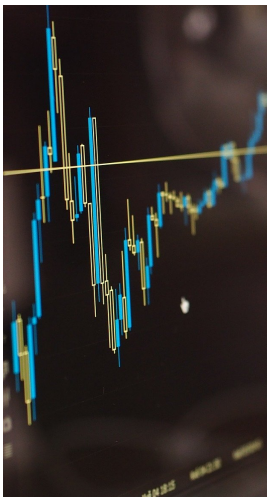


SOURCE: <https://pixabay.com/it/>

**Regression**, when possible, is a great method to learn how the distribution of a response  $y$  [or functionals of it], changes with covariates.

# Regression . . .

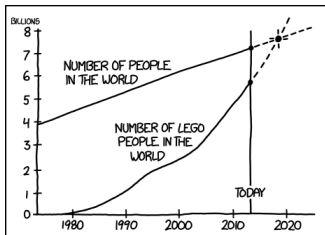
*“Statisticians are engaged in an exhausting but exhilarating struggle with the biggest challenge: **how to translate information into knowledge**” [S. Senn]*



SOURCE: <https://pixabay.com/it/>

**Regression**, when possible, is a great method to learn how the distribution of a response  $y$  [or functionals of it], changes with covariates.

**However**, going beyond regression for Gaussian responses [either from a frequentist or Bayesian perspective], introduces some issues.



BY 2019, HUMANS WILL BE OUTNUMBERED.

SOURCE: <https://xkcd.com/1281/>

# Bayesian probit regression

**Goal:** Given [conditionally] independent binary data  $y_1, \dots, y_n$  from a probit model  $(y_i | \beta) \sim \text{Bern}[\Phi(\mathbf{x}_i^T \beta)]$ ,  $i = 1, \dots, n$  with [in general] Gaussian prior  $\beta \sim N_p(\xi, \Omega)$  for  $\beta$ , **provide inference on the posterior**  $p(\beta | \mathbf{y})$ .

# Bayesian probit regression

**Goal:** Given [conditionally] independent binary data  $y_1, \dots, y_n$  from a probit model  $(y_i | \beta) \sim \text{Bern}[\Phi(\mathbf{x}_i^\top \beta)]$ ,  $i = 1, \dots, n$  with [in general] Gaussian prior  $\beta \sim N_p(\xi, \Omega)$  for  $\beta$ , **provide inference on the posterior**  $p(\beta | \mathbf{y})$ .

Applying [Bayes rule](#), the answer to the above question is

$$p(\beta | \mathbf{y}) = \frac{\phi_p(\beta - \xi; \Omega) \prod_{i=1}^n \Phi(\mathbf{x}_i^\top \beta)^{y_i} [1 - \Phi(\mathbf{x}_i^\top \beta)]^{1-y_i}}{\int_{\mathbb{R}^p} \phi_p(\beta - \xi; \Omega) \prod_{i=1}^n \Phi(\mathbf{x}_i^\top \beta)^{y_i} [1 - \Phi(\mathbf{x}_i^\top \beta)]^{1-y_i} d\beta}.$$

# Bayesian probit regression

**Goal:** Given [conditionally] independent binary data  $y_1, \dots, y_n$  from a probit model  $(y_i | \beta) \sim \text{Bern}[\Phi(\mathbf{x}_i^\top \beta)]$ ,  $i = 1, \dots, n$  with [in general] Gaussian prior  $\beta \sim N_p(\xi, \Omega)$  for  $\beta$ , **provide inference on the posterior**  $p(\beta | \mathbf{y})$ .

Applying **Bayes rule**, the answer to the above question is

$$p(\beta | \mathbf{y}) = \frac{\phi_p(\beta - \xi; \Omega) \prod_{i=1}^n \Phi(\mathbf{x}_i^\top \beta)^{y_i} [1 - \Phi(\mathbf{x}_i^\top \beta)]^{1-y_i}}{\int_{\mathbb{R}^p} \phi_p(\beta - \xi; \Omega) \prod_{i=1}^n \Phi(\mathbf{x}_i^\top \beta)^{y_i} [1 - \Phi(\mathbf{x}_i^\top \beta)]^{1-y_i} d\beta}.$$

**However**  $p(\beta | \mathbf{y})$  does not seem to belong to some known class of distributions and the normalizing constant apparently does not have an explicit form.

# Bayesian probit regression

**Goal:** Given [conditionally] independent binary data  $y_1, \dots, y_n$  from a probit model  $(y_i | \beta) \sim \text{Bern}[\Phi(\mathbf{x}_i^\top \beta)]$ ,  $i = 1, \dots, n$  with [in general] Gaussian prior  $\beta \sim N_p(\xi, \Omega)$  for  $\beta$ , **provide inference on the posterior**  $p(\beta | \mathbf{y})$ .

Applying **Bayes rule**, the answer to the above question is

$$p(\beta | \mathbf{y}) = \frac{\phi_p(\beta - \xi; \Omega) \prod_{i=1}^n \Phi(\mathbf{x}_i^\top \beta)^{y_i} [1 - \Phi(\mathbf{x}_i^\top \beta)]^{1-y_i}}{\int_{\mathbb{R}^p} \phi_p(\beta - \xi; \Omega) \prod_{i=1}^n \Phi(\mathbf{x}_i^\top \beta)^{y_i} [1 - \Phi(\mathbf{x}_i^\top \beta)]^{1-y_i} d\beta}.$$

**However**  $p(\beta | \mathbf{y})$  does not seem to belong to some known class of distributions and the normalizing constant apparently does not have an explicit form.

Arxiv:1407.0282v1 [stat.ML] 14 Jul 2014  
10.1145/2572488

## Leave Pima Indians Alone: Binary Regression as a Benchmark for Bayesian Computation

Nicolas Chopin and James Ridgway

Abstract. Whenever a new approach to perform Bayesian computation is introduced, a common practice is to showcase this approach on a binary regression model and datasets of moderate size. This paper discusses to which extent this practice is sound. It also reviews the current state of the art of Bayesian computation, using binary regression as a running example. Both sampling-based algorithms (importance sampling, MCMC and SMC) and fast approximations (Laplace, VB and EP) are covered. Extensive numerical results are provided, and are used to make recommendations to both end users and Bayesian computation experts. Implications for other problems (variable selection) and other models are also discussed.

Key words and phrases: Bayesian computation, expectation propagation.

**Solutions:** This has motivated several methods for **Bayesian inference in probit models**, covering **MCMC routines** [Metropolis–Hastings, Gibbs Sampling, Hamiltonian Monte Carlo] and **approximations of the posterior** [Laplace, Variational Bayes, Expectation Propagation].



# Unified skew-normal distribution

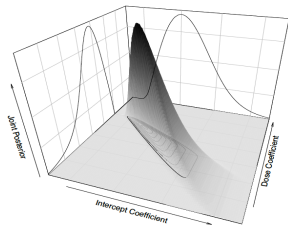
Arellano-Valle and Azzalini (2006), *Scandinavian Journal of Statistics*

Previous methods are still sub-optimal compared to cases in which the posterior belongs to a known and tractable class. **This could allow analytical posterior inference for Bayesian probit regression.** Indeed, the posterior is a SUN.

# Unified skew-normal distribution

Arellano-Valle and Azzalini (2006), *Scandinavian Journal of Statistics*

Previous methods are still sub-optimal compared to cases in which the posterior belongs to a known and tractable class. **This could allow analytical posterior inference for Bayesian probit regression.** Indeed, the posterior is a SUN.



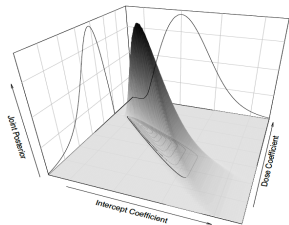
## Unified skew-normal random variable (SUN)

Generalizes the multivariate SN,  $\beta \sim \text{SN}_p(\xi, \Omega, \alpha)$  whose density  $2\phi_p(\beta - \xi; \Omega)\Phi[\alpha^\top \omega^{-1}(\beta - \xi)]$  is obtained by modifying a  $N_p(\xi, \Omega)$ , with the cdf of the  $N(0, 1)$  evaluated at  $\alpha^\top \omega^{-1}(\beta - \xi)$ , with  $\omega$  the diagonal matrix of standard deviations from  $\Omega$ . It unifies also other versions.

# Unified skew-normal distribution

Arellano-Valle and Azzalini (2006), *Scandinavian Journal of Statistics*

Previous methods are still sub-optimal compared to cases in which the posterior belongs to a known and tractable class. **This could allow analytical posterior inference for Bayesian probit regression.** Indeed, the posterior is a SUN.



## Unified skew-normal random variable (SUN)

Generalizes the multivariate SN,  $\beta \sim \text{SN}_p(\xi, \Omega, \alpha)$  whose density  $2\phi_p(\beta - \xi; \Omega)\Phi[\alpha^\top \omega^{-1}(\beta - \xi)]$  is obtained by modifying a  $N_p(\xi, \Omega)$ , with the cdf of the  $N(0, 1)$  evaluated at  $\alpha^\top \omega^{-1}(\beta - \xi)$ , with  $\omega$  the diagonal matrix of standard deviations from  $\Omega$ . It unifies also other versions.

More precisely, if  $\beta \sim \text{SUN}_{p,q}(\xi, \Omega, \Delta, \gamma, \Gamma)$ , with  $\xi \in \mathbb{R}^p$ ,  $\Delta \in \mathbb{R}^{p \times q}$ ,  $\gamma \in \mathbb{R}^q$  and  $\Omega^*$ —having block entries  $\Omega_{[11]}^* = \Gamma$ ,  $\Omega_{[22]}^* = \bar{\Omega}$  and  $\Omega_{[21]}^* = \Omega_{[12]}^{*\top} = \Delta$ —a full-rank correlation matrix, then the density is

$$\phi_p(\beta - \xi; \Omega) \frac{\Phi_q(\gamma + \Delta^\top \bar{\Omega}^{-1} \omega^{-1}(\beta - \xi); \Gamma - \Delta^\top \bar{\Omega}^{-1} \Delta)}{\Phi_q(\gamma; \Gamma)}, \quad (1)$$

# Unified skew-normal conjugacy in probit regression

Durante (2019), *Biometrika*

The posterior distribution  $p(\boldsymbol{\beta} \mid \mathbf{y})$  for the coefficients of a probit regression  $(y_i \mid \boldsymbol{\beta}) \sim \text{Bern}[\Phi(\mathbf{x}_i^\top \boldsymbol{\beta})]$ ,  $i = 1, \dots, n$ , coincides with a **unified skew-normal** (SUN) [Arellano-Valle and Azzalini, 2006], under Gaussian priors  $\boldsymbol{\beta} \sim N_p(\boldsymbol{\xi}, \boldsymbol{\Omega})$ .

# Unified skew-normal conjugacy in probit regression

Durante (2019), *Biometrika*

The posterior distribution  $p(\boldsymbol{\beta} \mid \mathbf{y})$  for the coefficients of a probit regression  $(y_i \mid \boldsymbol{\beta}) \sim \text{Bern}[\Phi(\mathbf{x}_i^\top \boldsymbol{\beta})]$ ,  $i = 1, \dots, n$ , coincides with a **unified skew-normal** (SUN) [Arellano-Valle and Azzalini, 2006], under Gaussian priors  $\boldsymbol{\beta} \sim N_p(\boldsymbol{\xi}, \boldsymbol{\Omega})$ .

**Main Theorem.** If  $(y_i \mid \boldsymbol{\beta}) \sim \text{Bern}[\Phi(\mathbf{x}_i^\top \boldsymbol{\beta})]$ ,  $i = 1, \dots, n$  and  $\boldsymbol{\beta} \sim N_p(\boldsymbol{\xi}, \boldsymbol{\Omega})$ :

$$(\boldsymbol{\beta} \mid \mathbf{y}) \sim \text{SUN}_{p,n}(\boldsymbol{\xi}, \boldsymbol{\Omega}, \bar{\boldsymbol{\Omega}} \boldsymbol{\omega} \mathbf{D}^\top \mathbf{s}^{-1}, \mathbf{s}^{-1} \mathbf{D} \boldsymbol{\xi}, \mathbf{s}^{-1} (\mathbf{D} \boldsymbol{\Omega} \mathbf{D}^\top + \mathbf{I}_n) \mathbf{s}^{-1}),$$

for every  $\mathbf{D} = \text{diag}(2y_1 - 1, \dots, 2y_n - 1) \mathbf{X} \in \mathbb{R}^{n \times p}$  and any  $n \times n$  positive diagonal matrix of standard deviations  $\mathbf{s} = [(\mathbf{D} \boldsymbol{\Omega} \mathbf{D}^\top + \mathbf{I}_n) \odot \mathbf{I}_n]^{1/2}$ .

# Unified skew-normal conjugacy in probit regression

Durante (2019), *Biometrika*

The **posterior distribution**  $p(\boldsymbol{\beta} \mid \mathbf{y})$  for the coefficients of a **probit regression**  $(y_i \mid \boldsymbol{\beta}) \sim \text{Bern}[\Phi(\mathbf{x}_i^\top \boldsymbol{\beta})]$ ,  $i = 1, \dots, n$ , coincides with a **unified skew-normal** (SUN) [Arellano-Valle and Azzalini, 2006], under **Gaussian priors**  $\boldsymbol{\beta} \sim N_p(\boldsymbol{\xi}, \boldsymbol{\Omega})$ .

**Main Theorem.** If  $(y_i \mid \boldsymbol{\beta}) \sim \text{Bern}[\Phi(\mathbf{x}_i^\top \boldsymbol{\beta})]$ ,  $i = 1, \dots, n$  and  $\boldsymbol{\beta} \sim N_p(\boldsymbol{\xi}, \boldsymbol{\Omega})$ :

$$(\boldsymbol{\beta} \mid \mathbf{y}) \sim \text{SUN}_{p,n}(\boldsymbol{\xi}, \boldsymbol{\Omega}, \bar{\boldsymbol{\Omega}} \boldsymbol{\omega} \mathbf{D}^\top \mathbf{s}^{-1}, \mathbf{s}^{-1} \mathbf{D} \boldsymbol{\xi}, \mathbf{s}^{-1} (\mathbf{D} \boldsymbol{\Omega} \mathbf{D}^\top + \mathbf{I}_n) \mathbf{s}^{-1}),$$

for every  $\mathbf{D} = \text{diag}(2y_1 - 1, \dots, 2y_n - 1) \mathbf{X} \in \mathbb{R}^{n \times p}$  and any  $n \times n$  positive diagonal matrix of standard deviations  $\mathbf{s} = [(\mathbf{D} \boldsymbol{\Omega} \mathbf{D}^\top + \mathbf{I}_n) \odot \mathbf{I}_n]^{1/2}$ .

**Sketch proof:** Note  $p(\boldsymbol{\beta} \mid \mathbf{y}) \propto \phi_p(\boldsymbol{\beta} - \boldsymbol{\xi}; \boldsymbol{\Omega}) \Phi_n(\mathbf{D} \boldsymbol{\beta}; \mathbf{I}_n)$  and that the kernel of  $\text{SUN}_{p,n}(\boldsymbol{\xi}, \boldsymbol{\Omega}, \boldsymbol{\Delta}, \boldsymbol{\gamma}, \boldsymbol{\Gamma})$  is  $\phi_p(\boldsymbol{\beta} - \boldsymbol{\xi}; \boldsymbol{\Omega}) \Phi_n(\boldsymbol{\gamma} + \boldsymbol{\Delta}^\top \bar{\boldsymbol{\Omega}}^{-1} \boldsymbol{\omega}^{-1} (\boldsymbol{\beta} - \boldsymbol{\xi}); \boldsymbol{\Gamma} - \boldsymbol{\Delta}^\top \bar{\boldsymbol{\Omega}}^{-1} \boldsymbol{\Delta})$ .

# Unified skew-normal conjugacy in probit regression

Durante (2019), *Biometrika*

The **posterior distribution**  $p(\boldsymbol{\beta} \mid \mathbf{y})$  for the coefficients of a **probit regression**  $(y_i \mid \boldsymbol{\beta}) \sim \text{Bern}[\Phi(\mathbf{x}_i^\top \boldsymbol{\beta})]$ ,  $i = 1, \dots, n$ , coincides with a **unified skew-normal** (SUN) [Arellano-Valle and Azzalini, 2006], under **Gaussian priors**  $\boldsymbol{\beta} \sim N_p(\boldsymbol{\xi}, \boldsymbol{\Omega})$ .

**Main Theorem.** If  $(y_i \mid \boldsymbol{\beta}) \sim \text{Bern}[\Phi(\mathbf{x}_i^\top \boldsymbol{\beta})]$ ,  $i = 1, \dots, n$  and  $\boldsymbol{\beta} \sim N_p(\boldsymbol{\xi}, \boldsymbol{\Omega})$ :

$$(\boldsymbol{\beta} \mid \mathbf{y}) \sim \text{SUN}_{p,n}(\boldsymbol{\xi}, \boldsymbol{\Omega}, \bar{\boldsymbol{\Omega}} \boldsymbol{\omega} \mathbf{D}^\top \mathbf{s}^{-1}, \mathbf{s}^{-1} \mathbf{D} \boldsymbol{\xi}, \mathbf{s}^{-1} (\mathbf{D} \boldsymbol{\Omega} \mathbf{D}^\top + \mathbf{I}_n) \mathbf{s}^{-1}),$$

for every  $\mathbf{D} = \text{diag}(2y_1 - 1, \dots, 2y_n - 1) \mathbf{X} \in \mathbb{R}^{n \times p}$  and any  $n \times n$  positive diagonal matrix of standard deviations  $\mathbf{s} = [(\mathbf{D} \boldsymbol{\Omega} \mathbf{D}^\top + \mathbf{I}_n) \odot \mathbf{I}_n]^{1/2}$ .

**Sketch proof:** Note  $p(\boldsymbol{\beta} \mid \mathbf{y}) \propto \phi_p(\boldsymbol{\beta} - \boldsymbol{\xi}; \boldsymbol{\Omega}) \Phi_n(\mathbf{D}\boldsymbol{\beta}; \mathbf{I}_n)$  and that the kernel of  $\text{SUN}_{p,n}(\boldsymbol{\xi}, \boldsymbol{\Omega}, \boldsymbol{\Delta}, \boldsymbol{\gamma}, \boldsymbol{\Gamma})$  is  $\phi_p(\boldsymbol{\beta} - \boldsymbol{\xi}; \boldsymbol{\Omega}) \Phi_n(\boldsymbol{\gamma} + \boldsymbol{\Delta}^\top \bar{\boldsymbol{\Omega}}^{-1} \boldsymbol{\omega}^{-1} (\boldsymbol{\beta} - \boldsymbol{\xi}); \boldsymbol{\Gamma} - \boldsymbol{\Delta}^\top \bar{\boldsymbol{\Omega}}^{-1} \boldsymbol{\Delta})$ .

**Remark:** **Whole SUN class is conjugate to probit.** Moreover, SUN has (i) **closure properties** [inference on  $(\boldsymbol{\beta}_j \mid \mathbf{y})$ ], (ii) **normalizing constant** fairly easy to compute [prediction and variable selection], (iii) simple **additive representation** [iid sampling], (iv) **explicit moment generating function** [posterior moments].

## Additive representation

To highlight the role of the hyperparameters  $\xi$  and  $\Omega$ , along with that of the data  $\mathbf{y}$  and  $\mathbf{X}$ , let us consider a [stochastic representation of the SUN posterior](#).

If  $(\beta | \mathbf{y}) \sim \text{SUN}_{p,n}(\xi, \Omega, \bar{\Omega}\omega\mathbf{D}^\top\mathbf{s}^{-1}, \mathbf{s}^{-1}\mathbf{D}\xi, \mathbf{s}^{-1}(\mathbf{D}\Omega\mathbf{D}^\top + \mathbf{I}_n)\mathbf{s}^{-1})$ , then

$$(\beta | \mathbf{y}) \stackrel{d}{=} \xi + \omega[\mathbf{V}_0 + \bar{\Omega}\omega\mathbf{D}^\top(\mathbf{D}\Omega\mathbf{D}^\top + \mathbf{I}_n)^{-1}\mathbf{s}\mathbf{V}_1], \quad (\mathbf{V}_0 \perp \mathbf{V}_1) \quad (2)$$

with  $\mathbf{V}_0 \sim N_p(\mathbf{0}, \bar{\Omega} - \bar{\Omega}\omega\mathbf{D}^\top(\mathbf{D}\Omega\mathbf{D}^\top + \mathbf{I}_n)^{-1}\mathbf{D}\omega\bar{\Omega})$ , and  $\mathbf{V}_1$  from an  $n$ -variate Gaussian  $N_n(\mathbf{0}, \mathbf{s}^{-1}(\mathbf{D}\Omega\mathbf{D}^\top + \mathbf{I}_n)\mathbf{s}^{-1})$  truncated below  $-\mathbf{s}^{-1}\mathbf{D}\xi$ .



# Additive representation

To highlight the role of the hyperparameters  $\xi$  and  $\Omega$ , along with that of the data  $\mathbf{y}$  and  $\mathbf{X}$ , let us consider a **stochastic representation of the SUN posterior**.

If  $(\beta | \mathbf{y}) \sim \text{SUN}_{p,n}(\xi, \Omega, \bar{\Omega}\omega\mathbf{D}^\top\mathbf{s}^{-1}, \mathbf{s}^{-1}\mathbf{D}\xi, \mathbf{s}^{-1}(\mathbf{D}\Omega\mathbf{D}^\top + \mathbf{I}_n)\mathbf{s}^{-1})$ , then

$$(\beta | \mathbf{y}) \stackrel{d}{=} \xi + \omega[\mathbf{V}_0 + \bar{\Omega}\omega\mathbf{D}^\top(\mathbf{D}\Omega\mathbf{D}^\top + \mathbf{I}_n)^{-1}\mathbf{s}\mathbf{V}_1], \quad (\mathbf{V}_0 \perp \mathbf{V}_1) \quad (2)$$

with  $\mathbf{V}_0 \sim N_p(\mathbf{0}, \bar{\Omega} - \bar{\Omega}\omega\mathbf{D}^\top(\mathbf{D}\Omega\mathbf{D}^\top + \mathbf{I}_n)^{-1}\mathbf{D}\omega\bar{\Omega})$ , and  $\mathbf{V}_1$  from an  $n$ -variate Gaussian  $N_n(\mathbf{0}, \mathbf{s}^{-1}(\mathbf{D}\Omega\mathbf{D}^\top + \mathbf{I}_n)\mathbf{s}^{-1})$  truncated below  $-\mathbf{s}^{-1}\mathbf{D}\xi$ .

**Comments:** The above representation provides some useful insights.

- $\xi$  has a main role on location, but has also an effect in controlling departures from normality both in terms of skewness and excess kurtosis.
- $\Omega$  has a main effect on scale and dependence, but contributes also to the shape in controlling the weight assigned to  $\mathbf{V}_1$ .
- Data in  $\mathbf{D}$  play more than a role in location, scale and departures from normality. If  $\mathbf{D} \approx \mathbf{0}$ ,  $\mathbf{V}_1$  has a negligible importance compared to  $\mathbf{V}_0$ .

**SUN is closed under marginalization, linear combinations and conditioning.** Adapting these results to the unified skew-normal in the previous theorem, the **marginal posteriors**  $(\beta_j | \mathbf{y})$ ,  $j = 1, \dots, p$ , still belong to the SUN family, and

$$\mathbb{E}(\beta | \mathbf{y}) = \xi + \Phi_n(\mathbf{s}^{-1} \mathbf{D} \xi; \mathbf{s}^{-1} (\mathbf{D} \Omega \mathbf{D}^\top + \mathbf{I}_n) \mathbf{s}^{-1})^{-1} \Omega \mathbf{D}^\top \mathbf{s}^{-1} \boldsymbol{\eta},$$

where  $\boldsymbol{\eta}$  is a simple function of the SUN parameters.

**SUN is closed under marginalization, linear combinations and conditioning.** Adapting these results to the unified skew-normal in the previous theorem, the **marginal posteriors**  $(\beta_j | \mathbf{y})$ ,  $j = 1, \dots, p$ , still belong to the SUN family, and

$$\mathbb{E}(\beta | \mathbf{y}) = \xi + \Phi_n(\mathbf{s}^{-1}\mathbf{D}\xi; \mathbf{s}^{-1}(\mathbf{D}\Omega\mathbf{D}^\top + \mathbf{I}_n)\mathbf{s}^{-1})^{-1}\Omega\mathbf{D}^\top\mathbf{s}^{-1}\eta,$$

where  $\eta$  is a simple function of the SUN parameters.

It is also possible to obtain closed-form expressions for **posterior predictive probabilities**  $\text{pr}(y_{\text{new}} = 1 | \mathbf{y}) = \int \Phi(\mathbf{x}_{\text{new}}^\top\beta)p(\beta | \mathbf{y})d\beta$  and the **marginal likelihood**  $\int p(\mathbf{y} | \mathcal{M}_k, \beta_{\mathcal{J}_k})p(\beta_{\mathcal{J}_k} | \mathcal{M}_k)d\beta_{\mathcal{J}_k}$  of a given model  $\mathcal{M}_k$ .

$$\text{pr}(y_{\text{new}} = 1 | \mathbf{y}) = \frac{\Phi_{n+1}(\mathbf{s}_{\text{new}}^{-1}\mathbf{D}_{\text{new}}\xi; \mathbf{s}_{\text{new}}^{-1}(\mathbf{D}_{\text{new}}\Omega\mathbf{D}_{\text{new}}^\top + \mathbf{I}_{n+1})\mathbf{s}_{\text{new}}^{-1})}{\Phi_n(\mathbf{s}^{-1}\mathbf{D}\xi; \mathbf{s}^{-1}(\mathbf{D}\Omega\mathbf{D}^\top + \mathbf{I}_n)\mathbf{s}^{-1})}.$$

The marginal likelihood is instead  $\Phi_n(\mathbf{s}_k^{-1}\mathbf{D}_k\xi_k; \mathbf{s}_k^{-1}(\mathbf{D}_k\Omega_k\mathbf{D}_k^\top + \mathbf{I}_n)\mathbf{s}_k^{-1})$ .

SUN is closed under marginalization, linear combinations and conditioning. Adapting these results to the unified skew-normal in the previous theorem, the marginal posteriors  $(\beta_j | \mathbf{y})$ ,  $j = 1, \dots, p$ , still belong to the SUN family, and

$$\mathbb{E}(\beta | \mathbf{y}) = \boldsymbol{\xi} + \Phi_n(\mathbf{s}^{-1} \mathbf{D} \boldsymbol{\xi}; \mathbf{s}^{-1} (\mathbf{D} \boldsymbol{\Omega} \mathbf{D}^\top + \mathbf{I}_n) \mathbf{s}^{-1})^{-1} \boldsymbol{\Omega} \mathbf{D}^\top \mathbf{s}^{-1} \boldsymbol{\eta},$$

where  $\boldsymbol{\eta}$  is a simple function of the SUN parameters.

It is also possible to obtain closed-form expressions for posterior predictive probabilities  $\text{pr}(y_{\text{new}} = 1 | \mathbf{y}) = \int \Phi(\mathbf{x}_{\text{new}}^\top \boldsymbol{\beta}) p(\boldsymbol{\beta} | \mathbf{y}) d\boldsymbol{\beta}$  and the marginal likelihood  $\int p(\mathbf{y} | \mathcal{M}_k, \boldsymbol{\beta}_{\mathcal{J}_k}) p(\boldsymbol{\beta}_{\mathcal{J}_k} | \mathcal{M}_k) d\boldsymbol{\beta}_{\mathcal{J}_k}$  of a given model  $\mathcal{M}_k$ .

$$\text{pr}(y_{\text{new}} = 1 | \mathbf{y}) = \frac{\Phi_{n+1}(\mathbf{s}_{\text{new}}^{-1} \mathbf{D}_{\text{new}} \boldsymbol{\xi}; \mathbf{s}_{\text{new}}^{-1} (\mathbf{D}_{\text{new}} \boldsymbol{\Omega} \mathbf{D}_{\text{new}}^\top + \mathbf{I}_{n+1}) \mathbf{s}_{\text{new}}^{-1})}{\Phi_n(\mathbf{s}^{-1} \mathbf{D} \boldsymbol{\xi}; \mathbf{s}^{-1} (\mathbf{D} \boldsymbol{\Omega} \mathbf{D}^\top + \mathbf{I}_n) \mathbf{s}^{-1})}.$$

The marginal likelihood is instead  $\Phi_n(\mathbf{s}_k^{-1} \mathbf{D}_k \boldsymbol{\xi}_k; \mathbf{s}_k^{-1} (\mathbf{D}_k \boldsymbol{\Omega}_k \mathbf{D}_k^\top + \mathbf{I}_n) \mathbf{s}_k^{-1})$ .

**Problem:** Inference requires sampling from  $n$ -variate truncated normals or evaluation of cumulative distribution functions  $\Phi_n(\cdot)$  of  $n$ -variate Gaussians.

# Closed-form filter for dynamic probit models

Fasano, Rebaudo, Durante, Petrone (2021), *Statistics and Computing*

**Goal:** Closed-form recursive expressions for predictive  $p(\beta_t | \mathbf{y}_{1:t-1})$ , filtering  $p(\beta_t | \mathbf{y}_{1:t})$  and smoothing  $p(\beta_{1:n} | \mathbf{y}_{1:n})$  distributions in the dynamic model

$$\begin{aligned}(y_t | \beta_t) &\sim \text{Bern}[\Phi(\mathbf{x}_t^T \beta_t)] \rightarrow p(y_t | \beta_t) = \Phi[(2y_t - 1)\mathbf{x}_t^T \beta_t] \\ \beta_t &= \mathbf{G}_t \beta_{t-1} + \varepsilon_t, \quad \varepsilon_t \sim N_p(\mathbf{0}, \mathbf{W}_t), \quad t = 1 \dots, n, \quad \beta_0 \sim N_p(\mathbf{a}_0, \mathbf{P}_0)\end{aligned}$$

# Closed-form filter for dynamic probit models

Fasano, Rebaudo, Durante, Petrone (2021), *Statistics and Computing*

**Goal:** Closed-form recursive expressions for predictive  $p(\beta_t | \mathbf{y}_{1:t-1})$ , filtering  $p(\beta_t | \mathbf{y}_{1:t})$  and smoothing  $p(\beta_{1:n} | \mathbf{y}_{1:n})$  distributions in the dynamic model

$$\begin{aligned}(y_t | \beta_t) &\sim \text{Bern}[\Phi(\mathbf{x}_t^T \beta_t)] \rightarrow p(y_t | \beta_t) = \Phi[(2y_t - 1)\mathbf{x}_t^T \beta_t] \\ \beta_t &= \mathbf{G}_t \beta_{t-1} + \varepsilon_t, \quad \varepsilon_t \sim N_p(\mathbf{0}, \mathbf{W}_t), \quad t = 1 \dots, n, \quad \beta_0 \sim N_p(\mathbf{a}_0, \mathbf{P}_0)\end{aligned}$$

**Hint:** Note that  $p(\beta_1 | y_1) \propto \phi_p(\beta_1 - \mathbf{G}_1 \mathbf{a}_0; \mathbf{G}_1 \mathbf{P}_0 \mathbf{G}_1^T + \mathbf{W}_1) \Phi[(2y_1 - 1)\mathbf{x}_1^T \beta_1]$ .

# Closed-form filter for dynamic probit models

Fasano, Rebaudo, Durante, Petrone (2021), *Statistics and Computing*

**Goal:** Closed-form recursive expressions for predictive  $p(\beta_t | \mathbf{y}_{1:t-1})$ , filtering  $p(\beta_t | \mathbf{y}_{1:t})$  and smoothing  $p(\beta_{1:n} | \mathbf{y}_{1:n})$  distributions in the dynamic model

$$(y_t | \beta_t) \sim \text{Bern}[\Phi(\mathbf{x}_t^\top \beta_t)] \rightarrow p(y_t | \beta_t) = \Phi[(2y_t - 1)\mathbf{x}_t^\top \beta_t]$$
$$\beta_t = \mathbf{G}_t \beta_{t-1} + \varepsilon_t, \quad \varepsilon_t \sim N_p(\mathbf{0}, \mathbf{W}_t), \quad t = 1 \dots, n, \quad \beta_0 \sim N_p(\mathbf{a}_0, \mathbf{P}_0)$$

**Hint:** Note that  $p(\beta_1 | y_1) \propto \phi_p(\beta_1 - \mathbf{G}_1 \mathbf{a}_0; \mathbf{G}_1 \mathbf{P}_0 \mathbf{G}_1^\top + \mathbf{W}_1) \Phi[(2y_1 - 1)\mathbf{x}_1^\top \beta_1]$ .

**Main Theorem** [closed-form filter for probit state-space models]

- 1 Filtering**  $[t - 1] \rightarrow$  **Predictive**  $[t]$ : If  $(\beta_{t-1} | \mathbf{y}_{1:t-1})$  is a  $\text{SUN}_{p,t-1}$  and  $\beta_t = \mathbf{G}_t \beta_{t-1} + \varepsilon_t$ , with  $\varepsilon_t \sim N_p(\mathbf{0}, \mathbf{W}_t)$ , then  $(\beta_t | \mathbf{y}_{1:t-1})$  is also a  $\text{SUN}_{p,t-1}$  with updated parameters [[closure under linear combinations](#)].
- 2 Predictive**  $[t] \rightarrow$  **Filtering**  $[t]$ : if  $(\beta_t | \mathbf{y}_{1:t-1})$  is  $\text{SUN}_{p,t-1}$  and  $p(y_t | \beta_t)$  is a probit likelihood, then  $p(\beta_t | \mathbf{y}_{1:t}) \propto p(\beta_t | \mathbf{y}_{1:t-1}) \Phi[(2y_t - 1)\mathbf{x}_t^\top \beta_t]$  is also  $\text{SUN}_{p,t}$  with updated parameters [[SUN-probit conjugacy](#); Durante, 2019].

# Closed-form filter for dynamic probit models

Fasano, Rebaudo, Durante, Petrone (2021), *Statistics and Computing*

**Goal:** Closed-form recursive expressions for predictive  $p(\beta_t | \mathbf{y}_{1:t-1})$ , filtering  $p(\beta_t | \mathbf{y}_{1:t})$  and smoothing  $p(\beta_{1:n} | \mathbf{y}_{1:n})$  distributions in the dynamic model

$$(y_t | \beta_t) \sim \text{Bern}[\Phi(\mathbf{x}_t^\top \beta_t)] \rightarrow p(y_t | \beta_t) = \Phi[(2y_t - 1)\mathbf{x}_t^\top \beta_t]$$
$$\beta_t = \mathbf{G}_t \beta_{t-1} + \varepsilon_t, \quad \varepsilon_t \sim N_p(\mathbf{0}, \mathbf{W}_t), \quad t = 1 \dots, n, \quad \beta_0 \sim N_p(\mathbf{a}_0, \mathbf{P}_0)$$

**Hint:** Note that  $p(\beta_1 | y_1) \propto \phi_p(\beta_1 - \mathbf{G}_1 \mathbf{a}_0; \mathbf{G}_1 \mathbf{P}_0 \mathbf{G}_1^\top + \mathbf{W}_1) \Phi[(2y_1 - 1)\mathbf{x}_1^\top \beta_1]$ .

**Main Theorem** [closed-form filter for probit state-space models]

- Filtering**  $[t - 1] \rightarrow$  **Predictive**  $[t]$ : If  $(\beta_{t-1} | \mathbf{y}_{1:t-1})$  is a  $\text{SUN}_{p,t-1}$  and  $\beta_t = \mathbf{G}_t \beta_{t-1} + \varepsilon_t$ , with  $\varepsilon_t \sim N_p(\mathbf{0}, \mathbf{W}_t)$ , then  $(\beta_t | \mathbf{y}_{1:t-1})$  is also a  $\text{SUN}_{p,t-1}$  with updated parameters [closure under linear combinations].
- Predictive**  $[t] \rightarrow$  **Filtering**  $[t]$ : if  $(\beta_t | \mathbf{y}_{1:t-1})$  is  $\text{SUN}_{p,t-1}$  and  $p(y_t | \beta_t)$  is a probit likelihood, then  $p(\beta_t | \mathbf{y}_{1:t}) \propto p(\beta_t | \mathbf{y}_{1:t-1}) \Phi[(2y_t - 1)\mathbf{x}_t^\top \beta_t]$  is also  $\text{SUN}_{p,t}$  with updated parameters [SUN-probit conjugacy; Durante, 2019].

Analog of the **Kalman filter** in the context of binary state-space models.



# SUN conjugacy in multinomial probit models

Fasano and Durante (2022), *Journal of Machine Learning Research*

Extension to  $L$  categories [based on **Gaussian latent utilities**  $u_{i1}, \dots, u_{iL}$ ].

- [Hausman and Wise, 1978].  $\text{pr}(y_i = l \mid \beta) = \text{pr}(u_{il} > u_{ik}, \forall k \neq l)$  with  $u_{il} = \mathbf{x}_{il}^T \beta + \varepsilon_{il}$ ,  $\varepsilon_i \sim N_L(\mathbf{0}, \Sigma)$  for each  $l = 1, \dots, L$  and  $i = 1, \dots, n$ .
- [Stern, 1992].  $\text{pr}(y_i = l \mid \beta) = \text{pr}(u_{il} > u_{ik}, \forall k \neq l)$  with  $u_{il} = \mathbf{x}_i^T \beta_l + \varepsilon_{il}$ ,  $\varepsilon_i \sim N_L(\mathbf{0}, \Sigma)$  for each  $l = 1, \dots, L$  and  $i = 1, \dots, n$ .
- [Tutz, 1991]. Based on a nested decision process relying on sequential binary decisions with probability  $\text{pr}(y_i = l \mid y_i > l - 1, \beta) = \Phi(\mathbf{x}_i^T \beta_l)$ .

# SUN conjugacy in multinomial probit models

Fasano and Durante (2022), *Journal of Machine Learning Research*

Extension to  $L$  categories [based on **Gaussian latent utilities**  $u_{i1}, \dots, u_{iL}$ ].

- [Hausman and Wise, 1978].  $\text{pr}(y_i = l \mid \beta) = \text{pr}(u_{il} > u_{ik}, \forall k \neq l)$  with  $u_{il} = \mathbf{x}_{il}^\top \beta + \varepsilon_{il}$ ,  $\varepsilon_i \sim N_L(\mathbf{0}, \Sigma)$  for each  $l = 1, \dots, L$  and  $i = 1, \dots, n$ .
- [Stern, 1992].  $\text{pr}(y_i = l \mid \beta) = \text{pr}(u_{il} > u_{ik}, \forall k \neq l)$  with  $u_{il} = \mathbf{x}_i^\top \beta_l + \varepsilon_{il}$ ,  $\varepsilon_i \sim N_L(\mathbf{0}, \Sigma)$  for each  $l = 1, \dots, L$  and  $i = 1, \dots, n$ .
- [Tutz, 1991]. Based on a nested decision process relying on sequential binary decisions with probability  $\text{pr}(y_i = l \mid y_i > l - 1, \beta) = \Phi(\mathbf{x}_i^\top \beta_l)$ .

**Goal:** Closed-form results for  $p(\beta \mid \mathbf{y})$ , when  $\beta$  has Gaussian or SUN prior.

# SUN conjugacy in multinomial probit models

Fasano and Durante (2022), *Journal of Machine Learning Research*

Extension to  $L$  categories [based on **Gaussian latent utilities**  $u_{i1}, \dots, u_{iL}$ ].

- [Hausman and Wise, 1978].  $\text{pr}(y_i = l \mid \beta) = \text{pr}(u_{il} > u_{ik}, \forall k \neq l)$  with  $u_{il} = \mathbf{x}_{il}^\top \beta + \varepsilon_{il}$ ,  $\varepsilon_i \sim N_L(\mathbf{0}, \Sigma)$  for each  $l = 1, \dots, L$  and  $i = 1, \dots, n$ .
- [Stern, 1992].  $\text{pr}(y_i = l \mid \beta) = \text{pr}(u_{il} > u_{ik}, \forall k \neq l)$  with  $u_{il} = \mathbf{x}_i^\top \beta_l + \varepsilon_{il}$ ,  $\varepsilon_i \sim N_L(\mathbf{0}, \Sigma)$  for each  $l = 1, \dots, L$  and  $i = 1, \dots, n$ .
- [Tutz, 1991]. Based on a nested decision process relying on sequential binary decisions with probability  $\text{pr}(y_i = l \mid y_i > l - 1, \beta) = \Phi(\mathbf{x}_i^\top \beta_l)$ .

**Goal:** Closed-form results for  $p(\beta \mid \mathbf{y})$ , when  $\beta$  has Gaussian or SUN prior.

**Hint:** The above models admit **likelihood** of the form  $p(\mathbf{y} \mid \beta) = \Phi_m(\bar{\mathbf{X}}\beta; \Lambda)$ , where  $\bar{\mathbf{X}}$  and  $\Lambda$  are suitable design and covariance matrices, respectively.

# SUN conjugacy in multinomial probit models

Fasano and Durante (2022), *Journal of Machine Learning Research*

Extension to  $L$  categories [based on **Gaussian latent utilities**  $u_{i1}, \dots, u_{iL}$ ].

- [Hausman and Wise, 1978].  $\text{pr}(y_i = l \mid \beta) = \text{pr}(u_{il} > u_{ik}, \forall k \neq l)$  with  $u_{il} = \mathbf{x}_{il}^\top \beta + \varepsilon_{il}$ ,  $\varepsilon_i \sim N_L(\mathbf{0}, \Sigma)$  for each  $l = 1, \dots, L$  and  $i = 1, \dots, n$ .
- [Stern, 1992].  $\text{pr}(y_i = l \mid \beta) = \text{pr}(u_{il} > u_{ik}, \forall k \neq l)$  with  $u_{il} = \mathbf{x}_i^\top \beta_l + \varepsilon_{il}$ ,  $\varepsilon_i \sim N_L(\mathbf{0}, \Sigma)$  for each  $l = 1, \dots, L$  and  $i = 1, \dots, n$ .
- [Tutz, 1991]. Based on a nested decision process relying on sequential binary decisions with probability  $\text{pr}(y_i = l \mid y_i > l - 1, \beta) = \Phi(\mathbf{x}_i^\top \beta_l)$ .

**Goal:** Closed-form results for  $p(\beta \mid \mathbf{y})$ , when  $\beta$  has Gaussian or SUN prior.

**Hint:** The above models admit **likelihood** of the form  $p(\mathbf{y} \mid \beta) = \Phi_m(\bar{\mathbf{X}}\beta; \Lambda)$ , where  $\bar{\mathbf{X}}$  and  $\Lambda$  are suitable design and covariance matrices, respectively.

**Main Theorem.** If  $p(\mathbf{y} \mid \beta) = \Phi_m(\bar{\mathbf{X}}\beta; \Lambda)$  and  $\beta$  has SUN prior (Gaussian is a special case), then  $(\beta \mid \mathbf{y}) \sim \text{SUN}_{q,m'}(\xi_{\text{POST}}, \Omega_{\text{POST}}, \Delta_{\text{POST}}, \gamma_{\text{POST}}, \Gamma_{\text{POST}})$ .

# SUN conjugacy in multinomial probit models

Fasano and Durante (2022), *Journal of Machine Learning Research*

Extension to  $L$  categories [based on **Gaussian latent utilities**  $u_{i1}, \dots, u_{iL}$ ].

- [Hausman and Wise, 1978].  $\text{pr}(y_i = l \mid \beta) = \text{pr}(u_{il} > u_{ik}, \forall k \neq l)$  with  $u_{il} = \mathbf{x}_{il}^\top \beta + \varepsilon_{il}$ ,  $\varepsilon_i \sim N_L(\mathbf{0}, \Sigma)$  for each  $l = 1, \dots, L$  and  $i = 1, \dots, n$ .
- [Stern, 1992].  $\text{pr}(y_i = l \mid \beta) = \text{pr}(u_{il} > u_{ik}, \forall k \neq l)$  with  $u_{il} = \mathbf{x}_i^\top \beta_l + \varepsilon_{il}$ ,  $\varepsilon_i \sim N_L(\mathbf{0}, \Sigma)$  for each  $l = 1, \dots, L$  and  $i = 1, \dots, n$ .
- [Tutz, 1991]. Based on a nested decision process relying on sequential binary decisions with probability  $\text{pr}(y_i = l \mid y_i > l - 1, \beta) = \Phi(\mathbf{x}_i^\top \beta_l)$ .

**Goal:** Closed-form results for  $p(\beta \mid \mathbf{y})$ , when  $\beta$  has Gaussian or SUN prior.

**Hint:** The above models admit **likelihood** of the form  $p(\mathbf{y} \mid \beta) = \Phi_m(\bar{\mathbf{X}}\beta; \Lambda)$ , where  $\bar{\mathbf{X}}$  and  $\Lambda$  are suitable design and covariance matrices, respectively.

**Main Theorem.** If  $p(\mathbf{y} \mid \beta) = \Phi_m(\bar{\mathbf{X}}\beta; \Lambda)$  and  $\beta$  has SUN prior (Gaussian is a special case), then  $(\beta \mid \mathbf{y}) \sim \text{SUN}_{q, m'}(\xi_{\text{POST}}, \Omega_{\text{POST}}, \Delta_{\text{POST}}, \gamma_{\text{POST}}, \Gamma_{\text{POST}})$ .

Leverage the **SUN properties** also for Bayesian inference in multinomial probits.

# Useful augmented–data representation

Albert and Chib (1993)

**Problem. Closed–form inference** under SUN posteriors requires to deal with **multivariate truncated normals** and **cdfs of multivariate Gaussians** whose dimension grows with the sample size  $n \rightarrow$  try to **approximate the posterior**.

# Useful augmented–data representation

Albert and Chib (1993)

**Problem.** Closed–form inference under SUN posteriors requires to deal with multivariate truncated normals and cdfs of multivariate Gaussians whose dimension grows with the sample size  $n \rightarrow$  try to **approximate the posterior**.

Bayesian probit regression model can also be expressed as

$$y_i = \mathbb{1}(z_i > 0), \text{ with } (z_i | \beta) \sim N(\mathbf{x}_i^T \beta, 1), i = 1, \dots, n, \text{ and } \beta \sim N_p(\mathbf{0}, \nu_p^2 \mathbf{I}_p).$$

Thus, we have a **dichotomized Gaussian linear regression on latent data**  $z_i$ .

# Useful augmented–data representation

Albert and Chib (1993)

**Problem.** Closed–form inference under SUN posteriors requires to deal with multivariate truncated normals and cdfs of multivariate Gaussians whose dimension grows with the sample size  $n \rightarrow$  try to approximate the posterior.

Bayesian probit regression model can also be expressed as

$$y_i = \mathbb{1}(z_i > 0), \text{ with } (z_i | \beta) \sim N(\mathbf{x}_i^T \beta, 1), i = 1, \dots, n, \text{ and } \beta \sim N_p(\mathbf{0}, \nu_p^2 \mathbf{I}_p).$$

Thus, we have a dichotomized Gaussian linear regression on latent data  $z_i$ .

This has been widely used in the development of MCMC and VB methods.

$$\begin{aligned} (\beta | \mathbf{z}, \mathbf{y}) &\sim N_p(\mathbf{V}\mathbf{X}^T\mathbf{z}, \mathbf{V}), \quad \mathbf{V} = (\nu_p^{-2}\mathbf{I}_p + \mathbf{X}^T\mathbf{X})^{-1}, \\ (z_i | \beta, \mathbf{z}_{-i}, \mathbf{y}) &\sim \begin{cases} \text{TN}[\mathbf{x}_i^T\beta, 1, (0, +\infty)], & \text{if } y_i = 1, \\ \text{TN}[\mathbf{x}_i^T\beta, 1, (-\infty, 0)], & \text{if } y_i = 0, \end{cases} \quad \text{for } i = 1, \dots, n, \end{aligned}$$

These full–conditionals allow implementation of Gibbs samplers [Albert and Chib, 1993] and mean–field VB with global and local variables [Consonni and Marin, 2007].



# Mean-field variational Bayes for probit models

**Goal:** Find a tractable approximation for the joint posterior density  $p(\beta, \mathbf{z} \mid \mathbf{y})$ , within the MF class of densities  $\mathcal{Q}_{\text{MF}} = \{q_{\text{MF}}(\beta, \mathbf{z}) : q_{\text{MF}}(\beta, \mathbf{z}) = q_{\text{MF}}(\beta)q_{\text{MF}}(\mathbf{z})\}$

# Mean-field variational Bayes for probit models

**Goal:** Find a tractable approximation for the joint posterior density  $p(\beta, \mathbf{z} \mid \mathbf{y})$ , within the MF class of densities  $\mathcal{Q}_{\text{MF}} = \{q_{\text{MF}}(\beta, \mathbf{z}) : q_{\text{MF}}(\beta, \mathbf{z}) = q_{\text{MF}}(\beta)q_{\text{MF}}(\mathbf{z})\}$

The optimal VB solution  $q_{\text{MF}}^*(\beta)q_{\text{MF}}^*(\mathbf{z})$  within this family minimizes

$$\text{KL}[q_{\text{MF}}(\beta, \mathbf{z}) \parallel p(\beta, \mathbf{z} \mid \mathbf{y})] = \mathbb{E}_{q_{\text{MF}}(\beta, \mathbf{z})}[\log q_{\text{MF}}(\beta, \mathbf{z})] - \mathbb{E}_{q_{\text{MF}}(\beta, \mathbf{z})}[\log p(\beta, \mathbf{z} \mid \mathbf{y})].$$

# Mean-field variational Bayes for probit models

**Goal:** Find a tractable approximation for the joint posterior density  $p(\beta, \mathbf{z} \mid \mathbf{y})$ , within the MF class of densities  $\mathcal{Q}_{\text{MF}} = \{q_{\text{MF}}(\beta, \mathbf{z}) : q_{\text{MF}}(\beta, \mathbf{z}) = q_{\text{MF}}(\beta)q_{\text{MF}}(\mathbf{z})\}$

The optimal VB solution  $q_{\text{MF}}^*(\beta)q_{\text{MF}}^*(\mathbf{z})$  within this family minimizes

$$\text{KL}[q_{\text{MF}}(\beta, \mathbf{z}) \parallel p(\beta, \mathbf{z} \mid \mathbf{y})] = \mathbb{E}_{q_{\text{MF}}(\beta, \mathbf{z})}[\log q_{\text{MF}}(\beta, \mathbf{z})] - \mathbb{E}_{q_{\text{MF}}(\beta, \mathbf{z})}[\log p(\beta, \mathbf{z} \mid \mathbf{y})].$$

In practice, we maximize  $\text{ELBO}[q_{\text{MF}}(\beta, \mathbf{z})] = -\text{KL}[q_{\text{MF}}(\beta, \mathbf{z}) \parallel p(\beta, \mathbf{z} \mid \mathbf{y})] + c$  via

$$q_{\text{MF}}^{(t)}(\beta) \propto \exp\{\mathbb{E}_{q_{\text{MF}}^{(t-1)}(\mathbf{z})} \log[p(\beta \mid \mathbf{z}, \mathbf{y})]\}, \quad q_{\text{MF}}^{(t)}(\mathbf{z}) \propto \exp\{\mathbb{E}_{q_{\text{MF}}^{(t)}(\beta)} \log[p(\mathbf{z} \mid \beta, \mathbf{y})]\},$$

that approximates  $p(\beta, \mathbf{z} \mid \mathbf{y})$  via  $q_{\text{MF}}^*(\beta) \prod_{i=1}^n q_{\text{MF}}^*(z_i)$ , where  $q_{\text{MF}}^*(\beta)$  is a Gaussian, while  $q_{\text{MF}}^*(z_i)$  are univariate truncated normals.

# Mean-field variational Bayes for probit models

**Goal:** Find a tractable approximation for the joint posterior density  $p(\beta, \mathbf{z} \mid \mathbf{y})$ , within the MF class of densities  $\mathcal{Q}_{\text{MF}} = \{q_{\text{MF}}(\beta, \mathbf{z}) : q_{\text{MF}}(\beta, \mathbf{z}) = q_{\text{MF}}(\beta)q_{\text{MF}}(\mathbf{z})\}$

The optimal VB solution  $q_{\text{MF}}^*(\beta)q_{\text{MF}}^*(\mathbf{z})$  within this family minimizes

$$\text{KL}[q_{\text{MF}}(\beta, \mathbf{z}) \parallel p(\beta, \mathbf{z} \mid \mathbf{y})] = \mathbb{E}_{q_{\text{MF}}(\beta, \mathbf{z})}[\log q_{\text{MF}}(\beta, \mathbf{z})] - \mathbb{E}_{q_{\text{MF}}(\beta, \mathbf{z})}[\log p(\beta, \mathbf{z} \mid \mathbf{y})].$$

In practice, we maximize ELBO  $[\log q_{\text{MF}}(\beta, \mathbf{z})] = -\text{KL}[q_{\text{MF}}(\beta, \mathbf{z}) \parallel p(\beta, \mathbf{z} \mid \mathbf{y})] + c$  via

$$q_{\text{MF}}^{(t)}(\beta) \propto \exp\{\mathbb{E}_{q_{\text{MF}}^{(t-1)}(\mathbf{z})} \log[p(\beta \mid \mathbf{z}, \mathbf{y})]\}, \quad q_{\text{MF}}^{(t)}(\mathbf{z}) \propto \exp\{\mathbb{E}_{q_{\text{MF}}^{(t)}(\beta)} \log[p(\mathbf{z} \mid \beta, \mathbf{y})]\},$$

that approximates  $p(\beta, \mathbf{z} \mid \mathbf{y})$  via  $q_{\text{MF}}^*(\beta) \prod_{i=1}^n q_{\text{MF}}^*(z_i)$ , where  $q_{\text{MF}}^*(\beta)$  is a Gaussian, while  $q_{\text{MF}}^*(z_i)$  are univariate truncated normals.

**However**, Fasano, Durante, Zanella (2022+) show that

# Mean-field variational Bayes for probit models

**Goal:** Find a tractable approximation for the joint posterior density  $p(\beta, \mathbf{z} \mid \mathbf{y})$ , within the MF class of densities  $\mathcal{Q}_{\text{MF}} = \{q_{\text{MF}}(\beta, \mathbf{z}) : q_{\text{MF}}(\beta, \mathbf{z}) = q_{\text{MF}}(\beta)q_{\text{MF}}(\mathbf{z})\}$

The optimal VB solution  $q_{\text{MF}}^*(\beta)q_{\text{MF}}^*(\mathbf{z})$  within this family minimizes

$$\text{KL}[q_{\text{MF}}(\beta, \mathbf{z}) \parallel p(\beta, \mathbf{z} \mid \mathbf{y})] = \mathbb{E}_{q_{\text{MF}}(\beta, \mathbf{z})}[\log q_{\text{MF}}(\beta, \mathbf{z})] - \mathbb{E}_{q_{\text{MF}}(\beta, \mathbf{z})}[\log p(\beta, \mathbf{z} \mid \mathbf{y})].$$

In practice, we maximize ELBO  $[\log q_{\text{MF}}(\beta, \mathbf{z})] = -\text{KL}[q_{\text{MF}}(\beta, \mathbf{z}) \parallel p(\beta, \mathbf{z} \mid \mathbf{y})] + c$  via

$$q_{\text{MF}}^{(t)}(\beta) \propto \exp\{\mathbb{E}_{q_{\text{MF}}^{(t-1)}(\mathbf{z})}[\log p(\beta \mid \mathbf{z}, \mathbf{y})]\}, \quad q_{\text{MF}}^{(t)}(\mathbf{z}) \propto \exp\{\mathbb{E}_{q_{\text{MF}}^{(t)}(\beta)}[\log p(\mathbf{z} \mid \beta, \mathbf{y})]\},$$

that approximates  $p(\beta, \mathbf{z} \mid \mathbf{y})$  via  $q_{\text{MF}}^*(\beta) \prod_{i=1}^n q_{\text{MF}}^*(z_i)$ , where  $q_{\text{MF}}^*(\beta)$  is a Gaussian, while  $q_{\text{MF}}^*(z_i)$  are univariate truncated normals.

**However**, Fasano, Durante, Zanella (2022+) show that

**Theorem:** Under simple assumptions,  $\liminf_{p \rightarrow \infty} \text{KL}[q_{\text{MF}}^*(\beta) \parallel p(\beta \mid \mathbf{y})] > 0$  almost surely (a.s.). Moreover,  $\nu_p^{-1} \|\mathbb{E}_{q_{\text{MF}}^*(\beta)}(\beta)\| \rightarrow 0$  (a.s.) as  $p \rightarrow \infty$ , where  $\|\cdot\|$  is the Euclidean norm. On the contrary,  $\nu_p^{-1} \|\mathbb{E}_{p(\beta \mid \mathbf{y})}(\beta)\| \rightarrow \text{const} \cdot \sqrt{n} > 0$  (a.s.) as  $p \rightarrow \infty$ , where **const** is a strictly positive constant.

# PFM-VB for probit models

Fasano, Durante, Zanella (2022+), *Biometrika*

**Solution:** Enlarge the class of approximating densities in a way that still allows simple optimization and inference. In particular, we consider the partially factorized family  $\mathcal{Q}_{\text{PFM}} = \{q_{\text{PFM}}(\beta, \mathbf{z}) : q_{\text{PFM}}(\beta, \mathbf{z}) = q_{\text{PFM}}(\beta | \mathbf{z}) \prod_{i=1}^n q_{\text{PFM}}(z_i)\}$ .

# PFM-VB for probit models

Fasano, Durante, Zanella (2022+), *Biometrika*

**Solution:** Enlarge the class of approximating densities in a way that still allows simple optimization and inference. In particular, we consider the partially factorized family  $\mathcal{Q}_{\text{PFM}} = \{q_{\text{PFM}}(\beta, \mathbf{z}) : q_{\text{PFM}}(\beta, \mathbf{z}) = q_{\text{PFM}}(\beta | \mathbf{z}) \prod_{i=1}^n q_{\text{PFM}}(z_i)\}$ .

**Motivation** for the use of  $\mathcal{Q}_{\text{PFM}}$ :  $q_{\text{MF}}^*(\beta, \mathbf{z}) = q_{\text{MF}}^*(\beta) \prod_{i=1}^n q_{\text{MF}}^*(z_i)$  belongs to  $\mathcal{Q}_{\text{PFM}}$ , and  $p(\beta, \mathbf{z} | \mathbf{y}) = p(\beta | \mathbf{z})p(\mathbf{z} | \mathbf{y})$  with  $p(\beta | \mathbf{z}) = \phi_p(\beta - \mathbf{V}\mathbf{X}^\top \mathbf{z}; \mathbf{V})$  and  $p(\mathbf{z} | \mathbf{y}) \propto \phi_n(\mathbf{z}; \mathbf{I}_n + \nu_p^2 \mathbf{X}\mathbf{X}^\top) \prod_{i=1}^n \mathbb{1}[(2y_i - 1)z_i > 0]$  [Holmes and Held, 2006].

# PFM-VB for probit models

Fasano, Durante, Zanella (2022+), *Biometrika*

**Solution:** Enlarge the class of approximating densities in a way that still allows simple optimization and inference. In particular, we consider the partially factorized family  $\mathcal{Q}_{\text{PFM}} = \{q_{\text{PFM}}(\beta, \mathbf{z}) : q_{\text{PFM}}(\beta, \mathbf{z}) = q_{\text{PFM}}(\beta | \mathbf{z}) \prod_{i=1}^n q_{\text{PFM}}(z_i)\}$ .

**Motivation** for the use of  $\mathcal{Q}_{\text{PFM}}$ :  $q_{\text{MF}}^*(\beta, \mathbf{z}) = q_{\text{MF}}^*(\beta) \prod_{i=1}^n q_{\text{MF}}^*(z_i)$  belongs to  $\mathcal{Q}_{\text{PFM}}$ , and  $p(\beta, \mathbf{z} | \mathbf{y}) = p(\beta | \mathbf{z})p(\mathbf{z} | \mathbf{y})$  with  $p(\beta | \mathbf{z}) = \phi_p(\beta - \mathbf{V}\mathbf{X}^T\mathbf{z}; \mathbf{V})$  and  $p(\mathbf{z} | \mathbf{y}) \propto \phi_n(\mathbf{z}; \mathbf{I}_n + \nu_p^2\mathbf{X}\mathbf{X}^T) \prod_{i=1}^n \mathbb{1}[(2y_i - 1)z_i > 0]$  [Holmes and Held, 2006].

**Proposition.** Let  $q_{\text{PFM}}^*(\beta, \mathbf{z})$  and  $q_{\text{MF}}^*(\beta, \mathbf{z})$  be the optimal approximations for  $p(\beta, \mathbf{z} | \mathbf{y})$ , under PFM-VB and MF-VB, respectively. Then

$$\text{KL}[q_{\text{PFM}}^*(\beta, \mathbf{z}) || p(\beta, \mathbf{z} | \mathbf{y})] \leq \text{KL}[q_{\text{MF}}^*(\beta, \mathbf{z}) || p(\beta, \mathbf{z} | \mathbf{y})].$$



# PFM-VB for probit models

Fasano, Durante, Zanella (2022+), *Biometrika*

**Solution:** Enlarge the class of approximating densities in a way that still allows simple optimization and inference. In particular, we consider the partially factorized family  $\mathcal{Q}_{\text{PFM}} = \{q_{\text{PFM}}(\beta, \mathbf{z}) : q_{\text{PFM}}(\beta, \mathbf{z}) = q_{\text{PFM}}(\beta | \mathbf{z}) \prod_{i=1}^n q_{\text{PFM}}(z_i)\}$ .

**Motivation** for the use of  $\mathcal{Q}_{\text{PFM}}$ :  $q_{\text{MF}}^*(\beta, \mathbf{z}) = q_{\text{MF}}^*(\beta) \prod_{i=1}^n q_{\text{MF}}^*(z_i)$  belongs to  $\mathcal{Q}_{\text{PFM}}$ , and  $p(\beta, \mathbf{z} | \mathbf{y}) = p(\beta | \mathbf{z})p(\mathbf{z} | \mathbf{y})$  with  $p(\beta | \mathbf{z}) = \phi_p(\beta - \mathbf{V}\mathbf{X}^T\mathbf{z}; \mathbf{V})$  and  $p(\mathbf{z} | \mathbf{y}) \propto \phi_n(\mathbf{z}; \mathbf{I}_n + \nu_p^2\mathbf{X}\mathbf{X}^T) \prod_{i=1}^n \mathbb{1}[(2y_i - 1)z_i > 0]$  [Holmes and Held, 2006].

**Proposition.** Let  $q_{\text{PFM}}^*(\beta, \mathbf{z})$  and  $q_{\text{MF}}^*(\beta, \mathbf{z})$  be the optimal approximations for  $p(\beta, \mathbf{z} | \mathbf{y})$ , under PFM-VB and MF-VB, respectively. Then

$$\text{KL}[q_{\text{PFM}}^*(\beta, \mathbf{z}) || p(\beta, \mathbf{z} | \mathbf{y})] \leq \text{KL}[q_{\text{MF}}^*(\beta, \mathbf{z}) || p(\beta, \mathbf{z} | \mathbf{y})].$$

**Main Theorem.** The optimal joint approximating density  $q_{\text{PFM}}^*(\beta, \mathbf{z})$  can be derived via a tractable **CAVI relying on simple closed-form expressions** and  $q_{\text{PFM}}^*(\beta) = \int_{\mathbb{R}^n} q_{\text{PFM}}^*(\beta | \mathbf{z}) \prod_{i=1}^n q_{\text{PFM}}^*(z_i) d\mathbf{z} = \mathbb{E}_{q_{\text{PFM}}^*(\mathbf{z})}[q_{\text{PFM}}^*(\beta | \mathbf{z})]$  of direct interest is the **density of a SUN, which crucially relies on a diagonal  $\mathbf{\Gamma} = \mathbf{I}_n$ .**

To be useful in practice,  $q_{\text{PFM}}^*(\beta, \mathbf{z})$  should be **simple to derive** and the density  $q_{\text{PFM}}^*(\beta) = \int_{\mathbb{R}^n} q_{\text{PFM}}^*(\beta | \mathbf{z}) \prod_{i=1}^n q_{\text{PFM}}^*(z_i) dz = \mathbb{E}_{q_{\text{PFM}}^*(\mathbf{z})}[q_{\text{PFM}}^*(\beta | \mathbf{z})]$  of direct interest should be **available in tractable form**.

To be useful in practice,  $q_{\text{PFM}}^*(\beta, \mathbf{z})$  should be **simple to derive** and the density  $q_{\text{PFM}}^*(\beta) = \int_{\mathbb{R}^n} q_{\text{PFM}}^*(\beta | \mathbf{z}) \prod_{i=1}^n q_{\text{PFM}}^*(z_i) dz = \mathbb{E}_{q_{\text{PFM}}^*(\mathbf{z})}[q_{\text{PFM}}^*(\beta | \mathbf{z})]$  of direct interest should be **available in tractable form**.

**Theorem:** Under the augmented probit model, the KL divergence between  $q_{\text{PFM}}(\beta, \mathbf{z}) \in \mathcal{Q}_{\text{PFM}}$  and  $p(\beta, \mathbf{z} | \mathbf{y})$  is minimized at  $q_{\text{PFM}}^*(\beta | \mathbf{z}) \prod_{i=1}^n q_{\text{PFM}}^*(z_i)$  with

$$q_{\text{PFM}}^*(\beta | \mathbf{z}) = p(\beta | \mathbf{z}) = \phi_{\rho}(\beta - \mathbf{V}\mathbf{X}^{\top}\mathbf{z}; \mathbf{V}), \quad \mathbf{V} = (\nu_{\rho}^{-2}\mathbf{I}_{\rho} + \mathbf{X}^{\top}\mathbf{X})^{-1},$$

$$q_{\text{PFM}}^*(z_i) = \frac{\phi(z_i - \mu_i^*; \sigma_i^{*2})}{\Phi[(2y_i - 1)\mu_i^*/\sigma_i^*]} \mathbb{1}[(2y_i - 1)z_i > 0], \quad \sigma_i^{*2} = (1 - \mathbf{x}_i^{\top}\mathbf{V}\mathbf{x}_i)^{-1},$$

where  $\boldsymbol{\mu}^* = (\mu_1^*, \dots, \mu_n^*)^{\top}$  solves  $\mu_i^* - \sigma_i^{*2}\mathbf{x}_i^{\top}\mathbf{V}\mathbf{X}_{-i}^{\top}\bar{\mathbf{z}}_{-i}^* = 0$ ,  $i = 1, \dots, n$ , with  $\mathbf{X}_{-i}$  the design matrix without the  $i$ th row, while  $\bar{\mathbf{z}}_{-i}^*$  is the  $(n-1) \times 1$  vector obtained by removing  $\bar{z}_i^* = \mu_i^* + (2y_i - 1)\sigma_i^*\phi(\mu_i^*/\sigma_i^*)\Phi[(2y_i - 1)\mu_i^*/\sigma_i^*]^{-1}$ ,  $i = 1, \dots, n$ , from the vector  $\bar{\mathbf{z}}^* = (\bar{z}_1^*, \dots, \bar{z}_n^*)^{\top}$ .

To be useful in practice,  $q_{\text{PFM}}^*(\beta, \mathbf{z})$  should be **simple to derive** and the density  $q_{\text{PFM}}^*(\beta) = \int_{\mathbb{R}^n} q_{\text{PFM}}^*(\beta | \mathbf{z}) \prod_{i=1}^n q_{\text{PFM}}^*(z_i) dz = \mathbb{E}_{q_{\text{PFM}}^*(\mathbf{z})}[q_{\text{PFM}}^*(\beta | \mathbf{z})]$  of direct interest should be **available in tractable form**.

**Theorem:** Under the augmented probit model, the KL divergence between  $q_{\text{PFM}}(\beta, \mathbf{z}) \in \mathcal{Q}_{\text{PFM}}$  and  $p(\beta, \mathbf{z} | \mathbf{y})$  is minimized at  $q_{\text{PFM}}^*(\beta | \mathbf{z}) \prod_{i=1}^n q_{\text{PFM}}^*(z_i)$  with

$$q_{\text{PFM}}^*(\beta | \mathbf{z}) = p(\beta | \mathbf{z}) = \phi_{\rho}(\beta - \mathbf{V}\mathbf{X}^T \mathbf{z}; \mathbf{V}), \quad \mathbf{V} = (\nu_{\rho}^{-2} \mathbf{I}_{\rho} + \mathbf{X}^T \mathbf{X})^{-1},$$

$$q_{\text{PFM}}^*(z_i) = \frac{\phi(z_i - \mu_i^*; \sigma_i^{*2})}{\Phi[(2y_i - 1)\mu_i^* / \sigma_i^*]} \mathbb{1}[(2y_i - 1)z_i > 0], \quad \sigma_i^{*2} = (1 - \mathbf{x}_i^T \mathbf{V} \mathbf{x}_i)^{-1},$$

where  $\boldsymbol{\mu}^* = (\mu_1^*, \dots, \mu_n^*)^T$  solves  $\mu_i^* - \sigma_i^{*2} \mathbf{x}_i^T \mathbf{V} \mathbf{X}_{-i}^T \bar{\mathbf{z}}_{-i}^* = 0$ ,  $i = 1, \dots, n$ , with  $\mathbf{X}_{-i}$  the design matrix without the  $i$ th row, while  $\bar{\mathbf{z}}_{-i}^*$  is the  $(n-1) \times 1$  vector obtained by removing  $\bar{\mathbf{z}}_i^* = \mu_i^* + (2y_i - 1)\sigma_i^* \phi(\mu_i^* / \sigma_i^*) \Phi[(2y_i - 1)\mu_i^* / \sigma_i^*]^{-1}$ ,  $i = 1, \dots, n$ , from the vector  $\bar{\mathbf{z}}^* = (\bar{\mathbf{z}}_1^*, \dots, \bar{\mathbf{z}}_n^*)^T$ .

The optimal parameters of the above densities can be obtained via a **simple CAVI algorithm** [at the same cost of MF-VB].

# Approximation quality and computational efficiency

Fasano, Durante, Zanella (2022+), *Biometrika*

The factorized form for  $q_{\text{PFM}}(\mathbf{z})$  leads to a SUN approximate density for  $\beta$ , with  $\Gamma = \mathbf{I}_n$ . This allows tractable inference at an  $\mathcal{O}(pn \cdot \min\{p, n\})$  cost.

# Approximation quality and computational efficiency

Fasano, Durante, Zanella (2022+), *Biometrika*

The factorized form for  $q_{\text{PFM}}(\mathbf{z})$  leads to a SUN approximate density for  $\beta$ , with  $\Gamma = \mathbf{I}_n$ . This allows tractable inference at an  $\mathcal{O}(pn \cdot \min\{p, n\})$  cost.

**Theorem** Under simple assumptions,  $\text{KL}[q_{\text{PFM}}^*(\beta) \parallel p(\beta \mid \mathbf{y})] \xrightarrow{p} 0$  as  $p \rightarrow \infty$   
[quality of the approximation]

# Approximation quality and computational efficiency

Fasano, Durante, Zanella (2022+), *Biometrika*

The factorized form for  $q_{\text{PFM}}(\mathbf{z})$  leads to a **SUN approximate density** for  $\beta$ , with  $\Gamma = \mathbf{I}_n$ . This allows tractable inference at an  $\mathcal{O}(pn \cdot \min\{p, n\})$  cost.

**Theorem** Under simple assumptions,  $\text{KL}[q_{\text{PFM}}^*(\beta) \parallel p(\beta | \mathbf{y})] \xrightarrow{p} 0$  as  $p \rightarrow \infty$   
[quality of the approximation]

**Corollary.** Let  $\text{pr}(y_{\text{NEW}} = 1 | \mathbf{y}) = \int \Phi(\mathbf{x}_{\text{NEW}}^T \beta) p(\beta | \mathbf{y}) d\beta$  be the exact posterior predictive probability for a new unit with predictors  $\mathbf{x}_{\text{NEW}} \in \mathbb{R}^p$ . Then, under simple assumptions,  $\sup_{\mathbf{x}_{\text{NEW}} \in \mathbb{R}^p} |\text{pr}_{\text{PFM}}(y_{\text{NEW}} = 1 | \mathbf{y}) - \text{pr}(y_{\text{NEW}} = 1 | \mathbf{y})| \xrightarrow{p} 0$  as  $p \rightarrow \infty$ . Instead,  $\liminf_{p \rightarrow \infty} \sup_{\mathbf{x}_{\text{NEW}} \in \mathbb{R}^p} |\text{pr}_{\text{MF}}(y_{\text{NEW}} = 1 | \mathbf{y}) - \text{pr}(y_{\text{NEW}} = 1 | \mathbf{y})| > 0$  almost surely as  $p \rightarrow \infty$  [quality of classification]

# Approximation quality and computational efficiency

Fasano, Durante, Zanella (2022+), *Biometrika*

The factorized form for  $q_{\text{PFM}}(\mathbf{z})$  leads to a **SUN approximate density** for  $\beta$ , with  $\Gamma = \mathbf{I}_n$ . This allows tractable inference at an  $\mathcal{O}(pn \cdot \min\{p, n\})$  cost.

**Theorem** Under simple assumptions,  $\text{KL}[q_{\text{PFM}}^*(\beta) \parallel p(\beta | \mathbf{y})] \xrightarrow{p} 0$  as  $p \rightarrow \infty$   
[quality of the approximation]

**Corollary.** Let  $\text{pr}(y_{\text{NEW}} = 1 | \mathbf{y}) = \int \Phi(\mathbf{x}_{\text{NEW}}^T \beta) p(\beta | \mathbf{y}) d\beta$  be the exact posterior predictive probability for a new unit with predictors  $\mathbf{x}_{\text{NEW}} \in \mathbb{R}^p$ . Then, under simple assumptions,  $\sup_{\mathbf{x}_{\text{NEW}} \in \mathbb{R}^p} |\text{pr}_{\text{PFM}}(y_{\text{NEW}} = 1 | \mathbf{y}) - \text{pr}(y_{\text{NEW}} = 1 | \mathbf{y})| \xrightarrow{p} 0$  as  $p \rightarrow \infty$ . Instead,  $\liminf_{p \rightarrow \infty} \sup_{\mathbf{x}_{\text{NEW}} \in \mathbb{R}^p} |\text{pr}_{\text{MF}}(y_{\text{NEW}} = 1 | \mathbf{y}) - \text{pr}(y_{\text{NEW}} = 1 | \mathbf{y})| > 0$  almost surely as  $p \rightarrow \infty$  [quality of classification]

**Theorem.** Let  $q_{\text{PFM}}^{(t)}(\beta) = \int_{\mathbb{R}^n} q_{\text{PFM}}^{(t)}(\beta | \mathbf{z}) \prod_{i=1}^n q_{\text{PFM}}^{(t)}(z_i) dz$  be the approximate density for  $\beta$  produced at iteration  $t$  by our CAVI. Then, under simple assumptions,  $\text{KL}[q_{\text{PFM}}^{(1)}(\beta) \parallel p(\beta | \mathbf{y})] \xrightarrow{p} 0$  as  $p \rightarrow \infty$  [computational efficiency]



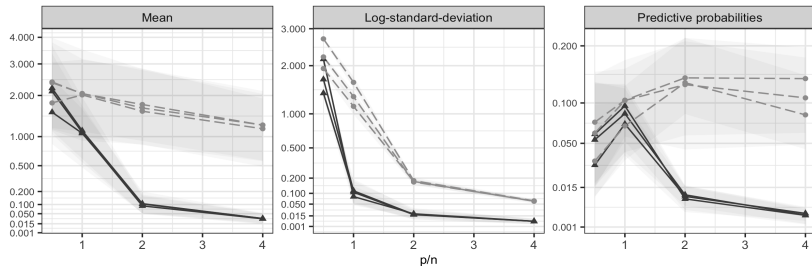
We evaluate accuracy in the approximation for three key functionals of the posterior distribution for  $\beta$ , by comparing MF-VB and PFM-VB approximations for these quantities with the STAN estimates at varying  $(p, n)$  settings.

**Simulation scenario:** data  $\mathbf{y}$  are simulated from probit regression with inputs  $x_{ij}$ ,  $[i = 1, \dots, n, j = 1, \dots, p]$  sampled from **independent standard normals** and coefficients  $\beta_j$   $[j = 1, \dots, p]$  simulated from uniforms in the range  $[-5, 5]$ .

# Simulation

We evaluate accuracy in the approximation for three key functionals of the posterior distribution for  $\beta$ , by comparing MF-VB and PFM-VB approximations for these quantities with the STAN estimates at varying  $(p, n)$  settings.

**Simulation scenario:** data  $\mathbf{y}$  are simulated from probit regression with inputs  $x_{ij}$ ,  $[i = 1, \dots, n, j = 1, \dots, p]$  sampled from **independent standard normals** and coefficients  $\beta_j$   $[j = 1, \dots, p]$  simulated from uniforms in the range  $[-5, 5]$ .



Empirical evidence is in line with theory and shows that our asymptotic results are visible also in finite-dimensional  $p > n$  settings.

# Alzheimers' application

**Large  $p$ , moderate  $n$  study** on presence–absence of Alzheimer as a function of demographic data, genotype and assay results. In this application  $n = 300$  and  $p = 9036$  [we include interactions]. We consider  $\beta \sim N_{9036}(\mathbf{0}, 25 \cdot \mathbf{I}_{9036})$ .

# Alzheimers' application

**Large  $p$ , moderate  $n$  study** on presence–absence of Alzheimer as a function of demographic data, genotype and assay results. In this application  $n = 300$  and  $p = 9036$  [we include interactions]. We consider  $\beta \sim N_{9036}(\mathbf{0}, 25 \cdot \mathbf{I}_{9036})$ .

**Computational performance.** Runtimes required for posterior inference

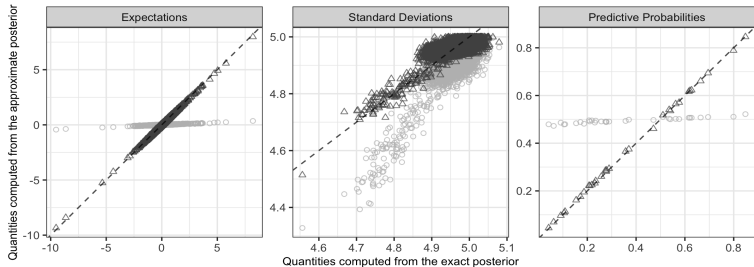
	STAN	EP	SUN	MF-VB	PFM-VB
Time [minutes]	> 360.00	> 360.00	92.27	0.04	0.04

# Alzheimers' application

**Large  $p$ , moderate  $n$  study** on presence–absence of Alzheimer as a function of demographic data, genotype and assay results. In this application  $n = 300$  and  $p = 9036$  [we include interactions]. We consider  $\beta \sim N_{9036}(\mathbf{0}, 25 \cdot \mathbf{I}_{9036})$ .

**Computational performance.** Runtimes required for posterior inference

	STAN	EP	SUN	MF-VB	PFM-VB
Time [minutes]	> 360.00	> 360.00	92.27	0.04	0.04



The models considered so far are special examples of a much **broader class of formulations** whose likelihood factorizes as

$$p(\mathbf{y} \mid \boldsymbol{\beta}) = p(\mathbf{y}_1 \mid \boldsymbol{\beta})p(\mathbf{y}_0 \mid \boldsymbol{\beta}) \propto \phi_{n_1}(\mathbf{y}_1 - \mathbf{X}_1\boldsymbol{\beta}; \boldsymbol{\Sigma}_1)\Phi_{n_0}(\mathbf{y}_0 + \mathbf{X}_0\boldsymbol{\beta}; \boldsymbol{\Sigma}_0). \quad (3)$$

**Examples:** probit, multivariate probit, multinomial probit, tobit, and others.

The models considered so far are special examples of a much **broader class of formulations** whose likelihood factorizes as

$$p(\mathbf{y} \mid \boldsymbol{\beta}) = p(\mathbf{y}_1 \mid \boldsymbol{\beta})p(\mathbf{y}_0 \mid \boldsymbol{\beta}) \propto \phi_{n_1}(\mathbf{y}_1 - \mathbf{X}_1\boldsymbol{\beta}; \boldsymbol{\Sigma}_1)\Phi_{n_0}(\mathbf{y}_0 + \mathbf{X}_0\boldsymbol{\beta}; \boldsymbol{\Sigma}_0). \quad (3)$$

**Examples:** probit, multivariate probit, multinomial probit, tobit, and others.

**Note:** Recalling the results in the previous slides, the above likelihood is actually the **kernel of a SUN density**  $\rightarrow$  **conjugacy** with SUN priors.

The models considered so far are special examples of a much **broader class of formulations** whose likelihood factorizes as

$$p(\mathbf{y} | \beta) = p(\mathbf{y}_1 | \beta)p(\mathbf{y}_0 | \beta) \propto \phi_{n_1}(\mathbf{y}_1 - \mathbf{X}_1\beta; \Sigma_1)\Phi_{n_0}(\mathbf{y}_0 + \mathbf{X}_0\beta; \Sigma_0). \quad (3)$$

**Examples:** probit, multivariate probit, multinomial probit, tobit, and others.

**Note:** Recalling the results in the previous slides, the above likelihood is actually the **kernel of a SUN density**  $\rightarrow$  **conjugacy** with SUN priors.

**Main Theorem.** If  $\beta \sim \text{SUN}_{p,q}(\xi, \Omega, \Delta, \gamma, \Gamma)$  — meaning that the prior density of  $\beta$  is  $p(\beta) \propto \phi_p(\beta - \xi; \Omega)\Phi_q(\gamma + \Delta^T \bar{\Omega}^{-1} \omega^{-1}(\beta - \xi); \Gamma - \Delta^T \bar{\Omega}^{-1} \Delta)$  — and  $p(\mathbf{y} | \beta)$  has likelihood (3), then

$$(\beta | \mathbf{y}) \sim \text{SUN}_{p,q+n_0}(\xi_{\text{POST}}, \Omega_{\text{POST}}, \Delta_{\text{POST}}, \gamma_{\text{POST}}, \Gamma_{\text{POST}}),$$

where  $\xi_{\text{POST}}, \Omega_{\text{POST}}, \Delta_{\text{POST}}, \gamma_{\text{POST}},$  and  $\Gamma_{\text{POST}}$  are simple analytical functions of  $\xi, \Omega, \Delta, \gamma, \Gamma$  and  $\mathbf{y}_1, \mathbf{X}_1, \Sigma_1, \mathbf{y}_0, \mathbf{X}_0, \Sigma_0$ .



The models considered so far are special examples of a much **broader class of formulations** whose likelihood factorizes as

$$p(\mathbf{y} \mid \beta) = p(\mathbf{y}_1 \mid \beta)p(\mathbf{y}_0 \mid \beta) \propto \phi_{n_1}(\mathbf{y}_1 - \mathbf{X}_1\beta; \Sigma_1)\Phi_{n_0}(\mathbf{y}_0 + \mathbf{X}_0\beta; \Sigma_0). \quad (3)$$

**Examples:** probit, multivariate probit, multinomial probit, tobit, and others.

**Note:** Recalling the results in the previous slides, the above likelihood is actually the **kernel of a SUN density**  $\rightarrow$  **conjugacy** with SUN priors.

**Main Theorem.** If  $\beta \sim \text{SUN}_{p,q}(\xi, \Omega, \Delta, \gamma, \Gamma)$  — meaning that the prior density of  $\beta$  is  $p(\beta) \propto \phi_p(\beta - \xi; \Omega)\Phi_q(\gamma + \Delta^T \bar{\Omega}^{-1} \omega^{-1}(\beta - \xi); \Gamma - \Delta^T \bar{\Omega}^{-1} \Delta)$  — and  $p(\mathbf{y} \mid \beta)$  has likelihood (3), then

$$(\beta \mid \mathbf{y}) \sim \text{SUN}_{p,q+n_0}(\xi_{\text{POST}}, \Omega_{\text{POST}}, \Delta_{\text{POST}}, \gamma_{\text{POST}}, \Gamma_{\text{POST}}),$$

where  $\xi_{\text{POST}}, \Omega_{\text{POST}}, \Delta_{\text{POST}}, \gamma_{\text{POST}},$  and  $\Gamma_{\text{POST}}$  are simple analytical functions of  $\xi, \Omega, \Delta, \gamma, \Gamma$  and  $\mathbf{y}_1, \mathbf{X}_1, \Sigma_1, \mathbf{y}_0, \mathbf{X}_0, \Sigma_0$ .

**Consequence:** All computational and inference methods previously developed can be applied to a broad class of routinely-implemented models.

# Other interesting results

Cao, Durante, Genton (2022+), *Journal of Computational and Graphical Statistics*

JOURNAL OF COMPUTATIONAL AND GRAPHICAL STATISTICS  
2022, VOL. 00, NO. 0, 1–12  
<https://doi.org/10.1080/10618600.2022.2036614>



Check for updates

## Scalable Computation of Predictive Probabilities in Probit Models with Gaussian Process Priors

Jian Cao<sup>a</sup>, Daniele Durante<sup>b</sup>, and Marc G. Genton<sup>a</sup>

<sup>a</sup>Statistics Program, King Abdullah University of Science and Technology, Thuwal, Saudi Arabia; <sup>b</sup>Department of Decision Sciences and Bocconi Institute for Data Science and Analytics, Bocconi University, Milano, Italy

### ABSTRACT

Predictive models for binary data are fundamental in various fields, and the growing complexity of modern applications has motivated several flexible specifications for modeling the relationship between the observed predictors and the binary responses. A widely-implemented solution is to express the probability parameter via a probit mapping of a Gaussian process indexed by predictors. However, unlike for continuous settings, there is a lack of closed-form results for predictive distributions in binary models with Gaussian process priors. Markov chain Monte Carlo methods and approximation strategies provide common solutions to this problem, but state-of-the-art algorithms are either computationally intractable or inaccurate in moderate-to-high dimensions. In this article, we aim to cover this gap by deriving closed-form expressions for the predictive probabilities in probit Gaussian processes that rely either on cumulative distribution functions of multivariate Gaussians or on functionals of multivariate truncated normals. To evaluate these quantities we develop novel scalable solutions based on tile-low-rank Monte Carlo methods for computing multivariate Gaussian probabilities, and on mean-field variational approximations of multivariate truncated normals. Closed-form expressions for the marginal likelihood and for the posterior distribution of the Gaussian process are also discussed. As shown in simulated and real-world empirical studies, the proposed methods scale to dimensions where state-of-the-art solutions are impractical.

### ARTICLE HISTORY

Received September 2020  
Revised November 2021

### KEYWORDS

Binary data; Gaussian process; Multivariate truncated normal; Probit model; Unified skew-normal; Variational Bayes

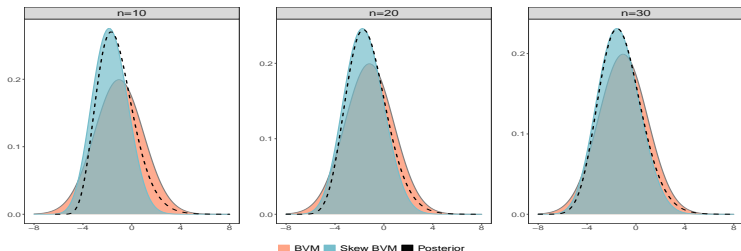
**Main result:** Derive closed-form expressions for the predictive probabilities in probit Gaussian processes that rely on ratios of cdfs of multivariate Gaussians and develop new scalable solutions based on **tile-low-rank Monte Carlo** methods and **separation-of-variables estimator** [Genz, 1992] for computing ratios of Gaussian cdfs with **theoretical accuracy guarantees**

# Bernstein–Von Mises theorem

**Bernstein–Von Mises theorem** [in short]: under regularity conditions, the **total variation distance** between the **posterior distribution** and a suitably–defined Gaussian distribution converges to **0** in probability.

# Bernstein–Von Mises theorem

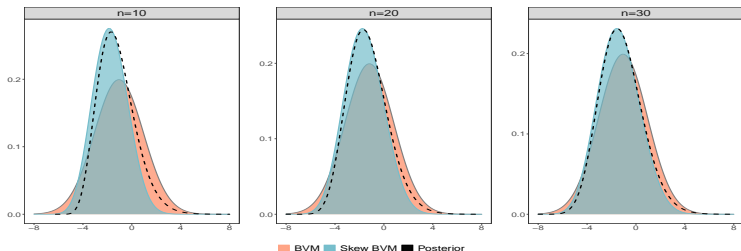
**Bernstein–Von Mises theorem** [in short]: under regularity conditions, the **total variation distance** between the **posterior distribution** and a suitably–defined **Gaussian distribution** converges to **0** in probability.



**However:** This limiting behavior may require a large sample size before becoming visible. In fact, **the posterior distribution is often skewed in practice.**

# Bernstein–Von Mises theorem

**Bernstein–Von Mises theorem** [in short]: under regularity conditions, the **total variation distance** between the **posterior distribution** and a suitably–defined **Gaussian distribution** converges to **0** in probability.



**However:** This limiting behavior may require a large sample size before becoming visible. In fact, **the posterior distribution is often skewed in practice.**

**Conjecture:** Adopting as limiting law a **skewed generalization of the Gaussian distribution**, we might obtain substantially more accurate/stronger results.

# Skewed Bernstein–Von Mises theorem

Pozza, Durante, Szabo (2022+), *soon online*

Let  $\{\mathbf{y}_i\}_{i=1}^n$  be a sequence of independent random variables with probability measure  $P_{\theta_0}^{(n)} \in \{P_{\theta}^{(n)}, \theta \in \Theta \subseteq \mathbb{R}^p\}$ . Moreover, let  $\ell(\theta)$  be the log-likelihood and  $\ell^{(1)} = [\ell_r^{(1)}]$ ,  $\ell^{(2)} = [\ell_{rs}^{(2)}]$ ,  $\ell^{(3)} = [\ell_{rst}^{(3)}]$  its first three derivatives at  $\theta_0$ .

# Skewed Bernstein–Von Mises theorem

Pozza, Durante, Szabo (2022+), *soon online*

Let  $\{\mathbf{y}_i\}_{i=1}^n$  be a sequence of independent random variables with probability measure  $P_{\theta_0}^{(n)} \in \{P_{\theta}^{(n)}, \theta \in \Theta \subseteq \mathbb{R}^p\}$ . Moreover, let  $\ell(\theta)$  be the log-likelihood and  $\ell^{(1)} = [\ell_r^{(1)}]$ ,  $\ell^{(2)} = [\ell_{rs}^{(2)}]$ ,  $\ell^{(3)} = [\ell_{rst}^{(3)}]$  its first three derivatives at  $\theta_0$ .

**Theorem:** Under regularity conditions on the log-likelihood ratio and its derivatives, if the map  $\theta \rightarrow P_{\theta}^{(n)}$  is one-to-one,  $\theta_0$  is an inner point of  $\Theta$  and the prior measure  $P(\theta)$  is absolutely continuous with bounded and positive density in a neighborhood of  $\theta_0$ , then

$$\| P(\cdot | \mathbf{y}^{(n)}) - P_{se}(\cdot) \|_{TV} = O_p(\{\log n\}^{p/2+3}/n)$$

where  $P_{se}(\mathbb{A}) = \int_{\mathbb{A}} p_{se}(\bar{\theta}) d\bar{\theta}$  for  $\mathbb{A} \subset \mathbb{R}^p$ ,  $\bar{\theta} = \sqrt{n}(\theta - \theta_0)$  and  $p_{se}(\bar{\theta})$  is the density of a suitably-defined skew-symmetric distribution [Azzalini & Regoli, 2012]. Specifically,  $p_{se}(\bar{\theta}) = 2\phi_p(\bar{\theta}; \xi_n, \Omega_n)\Phi\{\alpha_n(\bar{\theta})\}$ , where  $\alpha_n(\cdot)$  is an odd function.

# Skewed Bernstein–Von Mises theorem

Pozza, Durante, Szabo (2022+), *soon online*

Let  $\{\mathbf{y}_i\}_{i=1}^n$  be a sequence of independent random variables with probability measure  $P_{\theta_0}^{(n)} \in \{P_{\theta}^{(n)}, \theta \in \Theta \subseteq \mathbb{R}^p\}$ . Moreover, let  $\ell(\theta)$  be the log-likelihood and  $\ell^{(1)} = [\ell_r^{(1)}]$ ,  $\ell^{(2)} = [\ell_{rs}^{(2)}]$ ,  $\ell^{(3)} = [\ell_{rst}^{(3)}]$  its first three derivatives at  $\theta_0$ .

**Theorem:** Under regularity conditions on the log-likelihood ratio and its derivatives, if the map  $\theta \rightarrow P_{\theta}^{(n)}$  is one-to-one,  $\theta_0$  is an inner point of  $\Theta$  and the prior measure  $P(\theta)$  is absolutely continuous with bounded and positive density in a neighborhood of  $\theta_0$ , then

$$\|P(\cdot | \mathbf{y}^{(n)}) - P_{se}(\cdot)\|_{TV} = O_p(\{\log n\}^{p/2+3}/n)$$

where  $P_{se}(\mathbb{A}) = \int_{\mathbb{A}} p_{se}(\bar{\theta}) d\bar{\theta}$  for  $\mathbb{A} \subset \mathbb{R}^p$ ,  $\bar{\theta} = \sqrt{n}(\theta - \theta_0)$  and  $p_{se}(\bar{\theta})$  is the density of a suitably-defined skew-symmetric distribution [Azzalini & Regoli, 2012]. Specifically,  $p_{se}(\bar{\theta}) = 2\phi_p(\bar{\theta}; \xi_n, \Omega_n)\Phi\{\alpha_n(\bar{\theta})\}$ , where  $\alpha_n(\cdot)$  is an odd function.

**Remark:** In the above theorem, the quantities  $\xi_n$ ,  $\Omega_n$  and  $\alpha_n(\cdot)$  are simple analytical functions of  $\ell^{(1)} = [\ell_r^{(1)}]$ ,  $\ell^{(2)} = [\ell_{rs}^{(2)}]$ ,  $\ell^{(3)} = [\ell_{rst}^{(3)}]$  and the prior.



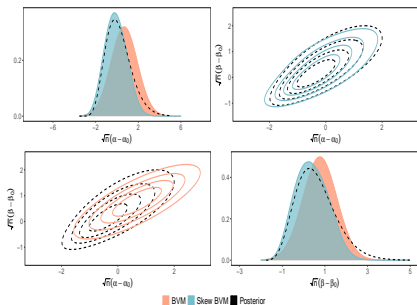
# Skew-modal approximation

Pozza, Durante, Szabo (2022+), *soon online*



# Skew-modal approximation

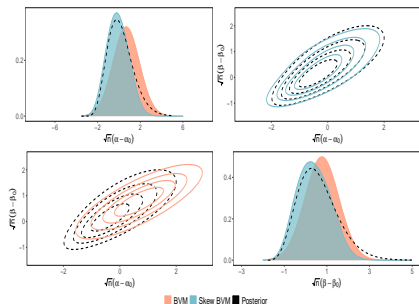
Pozza, Durante, Szabo (2022+), *soon online*



**Simulation** with  $n = 15$ ,  $y_i \stackrel{iid}{\sim} \text{Ga}(\alpha, \beta)$ ,  
 $\alpha \sim \text{Ga}(2)$  and  $\beta \sim \text{Ga}(2)$ .

# Skew-modal approximation

Pozza, Durante, Szabo (2022+), *soon online*



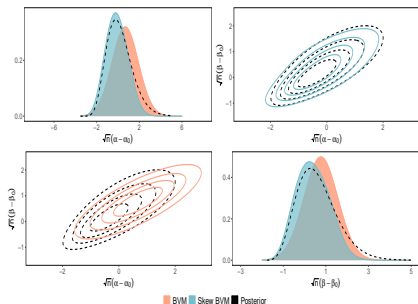
**Simulation** with  $n = 15$ ,  $y_i \stackrel{iid}{\sim} \text{Ga}(\alpha, \beta)$ ,  $\alpha \sim \text{Ga}(2)$  and  $\beta \sim \text{Ga}(2)$ .

**Comment.** We improve the approximation accuracy relative to classical BvM. However, both approximations require  $\theta_0$ , which is **not available in practice**.

**Solution.** Modal approximation based on a **skew-symmetric density** rather than a Gaussian one [recall Laplace approximation]

# Skew-modal approximation

Pozza, Durante, Szabo (2022+), soon online



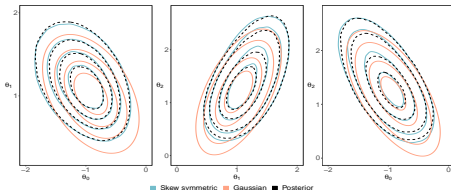
**Simulation** with  $n = 15$ ,  $y_i \stackrel{iid}{\sim} \text{Ga}(\alpha, \beta)$ ,  $\alpha \sim \text{Ga}(2)$  and  $\beta \sim \text{Ga}(2)$ .

**Comment.** We improve the approximation accuracy relative to classical BvM. However, both approximations require  $\theta_0$ , which is **not available in practice**.

**Solution.** Modal approximation based on a **skew-symmetric density** rather than a Gaussian one [recall Laplace approximation]

## Skew-modal approximation [provably

more accurate than Laplace]: Let  $\tilde{\ell}$  denote the log-posterior at its MAP  $\tilde{\theta}$ , then we approximate  $p(\theta | \mathbf{y}^{(n)})$  via  $2\phi_p(\theta; \tilde{\theta}, \tilde{\Omega})\Phi\{\tilde{\alpha}(\theta)\}$  where  $\tilde{\Omega}$  and  $\tilde{\alpha}(\cdot)$  are simple functions of  $\tilde{\ell}^{(2)}$ ,  $\tilde{\ell}^{(3)}$  and  $\tilde{\theta}$ .



**Main message:** Skew-normals and related families [Azzalini & co-authors] play a key role in Bayesian inference, which has been partially overlooked to date [Exception: Liseo & co-authors]. The advancements presented open new avenues for improved posterior inference via novel closed-form expressions, new Monte Carlo methods, and more accurate and scalable approximations.

**Main message:** Skew-normals and related families [Azzalini & co-authors] play a key role in Bayesian inference, which has been partially overlooked to date [Exception: Liseo & co-authors]. The advancements presented open new avenues for improved posterior inference via novel closed-form expressions, new Monte Carlo methods, and more accurate and scalable approximations.

The above results also motivate **further extensions**.

- Further improve the skew-modal approximation in terms of accuracy
- Explore conjugacy in broader classes [of models and skewed prior]
- Explore more complex models building on such representations; i.e. BART

**Main message:** Skew-normals and related families [Azzalini & co-authors] play a key role in Bayesian inference, which has been partially overlooked to date [Exception: Liseo & co-authors]. The advancements presented open new avenues for improved posterior inference via novel closed-form expressions, new Monte Carlo methods, and more accurate and scalable approximations.

The above results also motivate **further extensions**.

- Further improve the skew-modal approximation in terms of accuracy
- Explore conjugacy in broader classes [of models and skewed prior]
- Explore more complex models building on such representations; i.e. BART

Thank you for the attention!

<https://danieledurante.github.io/web/>

<https://github.com/danieledurante>

I would like to acknowledge the support from MIUR-PRIN 2017 project SELECT [20177BRJXS] in the preparation of some of these works.