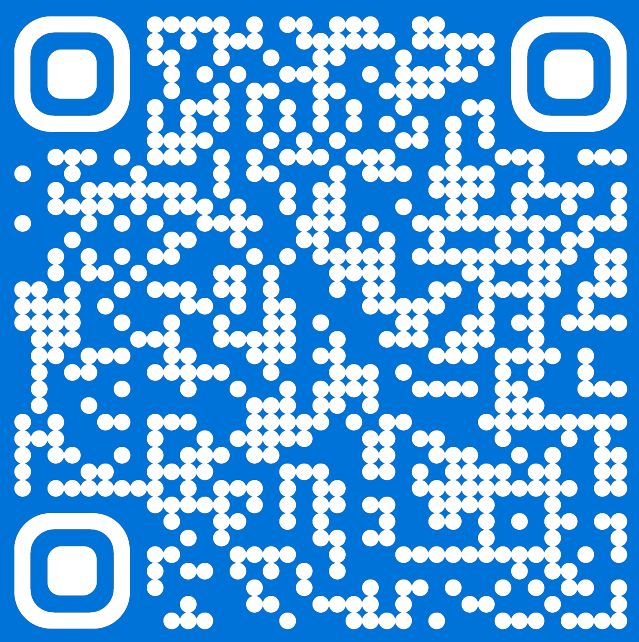


# We perform exact inference and hyperparameter optimisation in Bayesian linear models with millions of parameters



Take a picture to see the full paper.

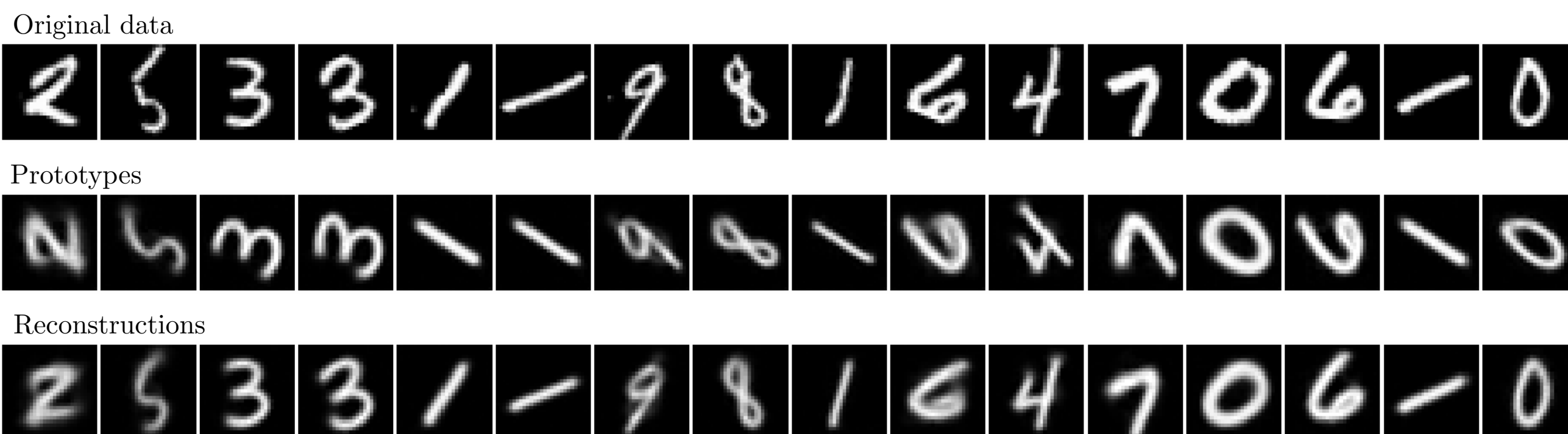
## Sampling-based inference for large linear models, with application to linearised Laplace



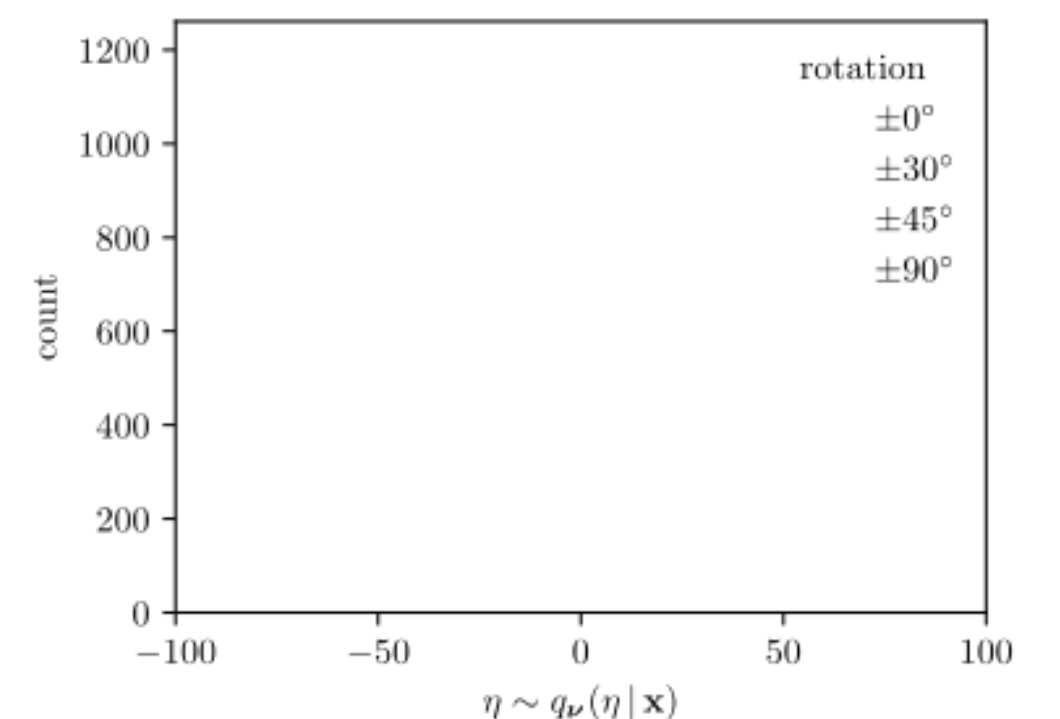
Javier Antorán\*, Shreyas Padhy,\* Riccardo Barbano, Eric Nalisnick, David Janz, and José Miguel Hernández-Lobato



We learn representations of MNIST digits that are invariant to rotations and successfully reconstruct the original digits:



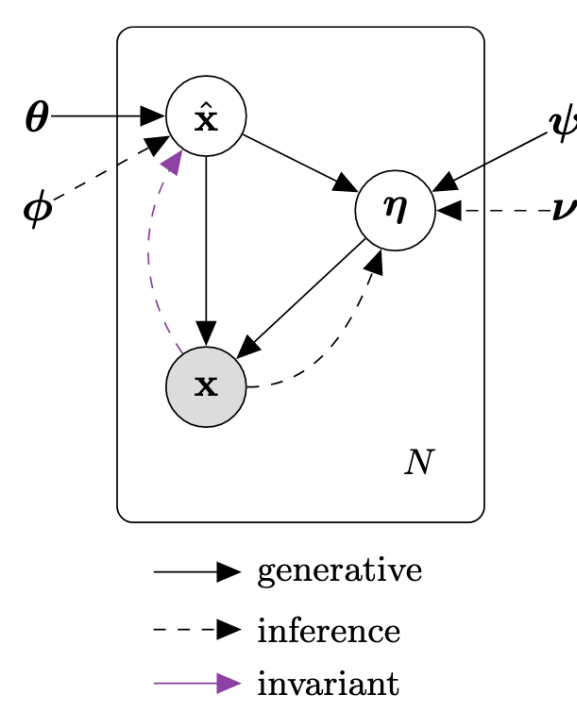
As digits are rotated more, we learn to predict larger angles:



Our model:

$$\begin{aligned} \hat{\mathbf{x}} &\sim p_{\theta}(\hat{\mathbf{x}}), \\ \boldsymbol{\eta} &\sim p_{\psi}(\boldsymbol{\eta} | \hat{\mathbf{x}}), \\ \mathbf{x} &\sim p(\mathbf{x} | \boldsymbol{\eta}, \hat{\mathbf{x}}), \end{aligned}$$

(1)  
(2)  
(3)



Affine transformation parameterization:

$$T_{\boldsymbol{\eta}}(\hat{\mathbf{x}}) = T_{\boldsymbol{\eta}} \cdot \hat{\mathbf{x}}, \quad T_{\boldsymbol{\eta}} = \exp \left( \sum_i \eta_i G_i \right) \quad (4)$$

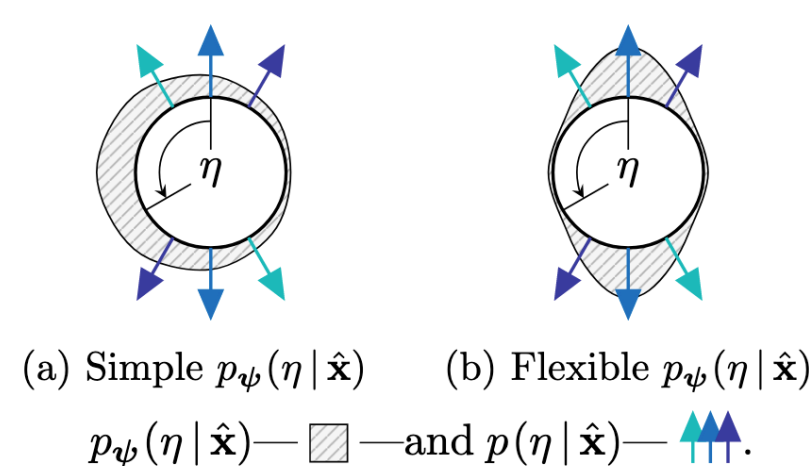
Our training objective:

$$\log p(\mathbf{x}) = \log \iint p(\mathbf{x}, \hat{\mathbf{x}}, \boldsymbol{\eta}) d\boldsymbol{\eta} d\hat{\mathbf{x}} \quad (5)$$

$$= \log \mathbb{E}_{q_{\nu}(\boldsymbol{\eta} | \mathbf{x}) q_{\phi}(\hat{\mathbf{x}} | \mathbf{x})} \left[ \frac{p(\mathbf{x} | \hat{\mathbf{x}}, \boldsymbol{\eta}) p_{\psi}(\boldsymbol{\eta} | \hat{\mathbf{x}}) p_{\theta}(\hat{\mathbf{x}})}{q_{\nu}(\boldsymbol{\eta} | \mathbf{x}) q_{\phi}(\hat{\mathbf{x}} | \mathbf{x})} \right] \quad (6)$$

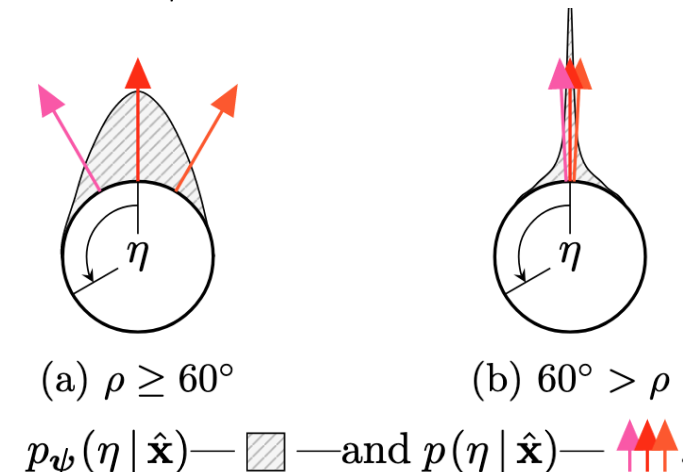
$$\geq \mathbb{E}_{q_{\nu} q_{\phi}} [\log p(\mathbf{x} | \hat{\mathbf{x}}, \boldsymbol{\eta})] - \mathbb{E}_{q_{\phi}} [D_{\text{KL}}(q_{\nu} || p_{\psi})] - D_{\text{KL}}(q_{\phi} || p_{\theta}) \equiv -\mathcal{L}(\boldsymbol{\theta}, \boldsymbol{\psi}, \boldsymbol{\phi}, \boldsymbol{\nu}) \quad (7)$$

Conjecture 1:  $p_{\psi}(\boldsymbol{\eta} | \hat{\mathbf{x}})$  must be sufficiently flexible.



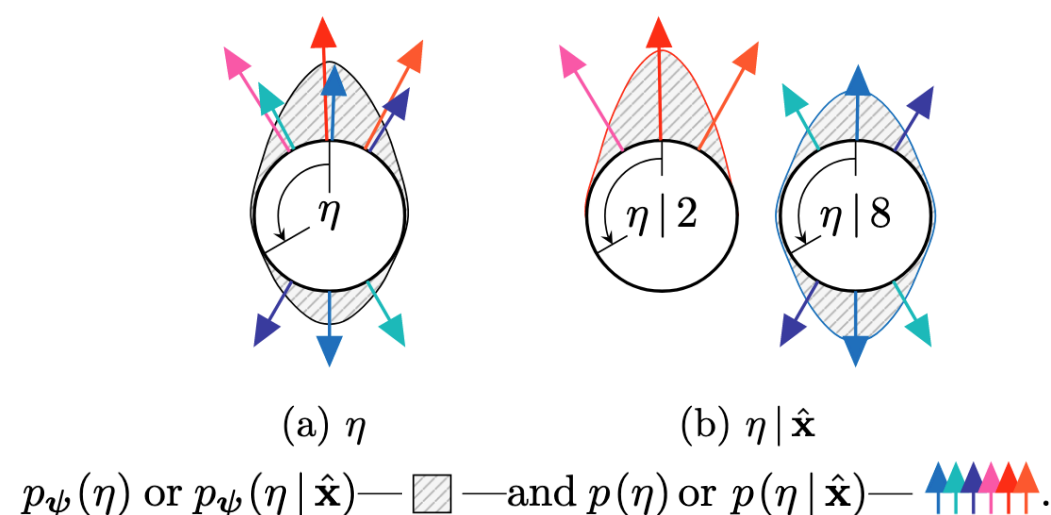
$\mathbf{x}$	$\hat{\mathbf{x}}$	$p(\boldsymbol{\eta}   \mathbf{x}, \hat{\mathbf{x}})$
8	8	$0.5 \cdot \delta(\boldsymbol{\eta} - 0^\circ) + 0.5 \cdot \delta(\boldsymbol{\eta} - 180^\circ)$
8	8	$0.5 \cdot \delta(\boldsymbol{\eta} - 30^\circ) + 0.5 \cdot \delta(\boldsymbol{\eta} + 150^\circ)$
8	8	$0.5 \cdot \delta(\boldsymbol{\eta} + 30^\circ) + 0.5 \cdot \delta(\boldsymbol{\eta} - 150^\circ)$

Conjecture 2:  $q_{\phi}(\hat{\mathbf{x}} | \mathbf{x})$  must be *fully* invariant w.r.t  $\boldsymbol{\eta}$ .



$\mathbf{x}$	$\hat{\mathbf{x}}$	$\rho \geq 60^\circ$		$60^\circ > \rho$	
		$p(\boldsymbol{\eta}   \mathbf{x}, \hat{\mathbf{x}})$	$\hat{\mathbf{x}}$	$p(\boldsymbol{\eta}   \mathbf{x}, \hat{\mathbf{x}})$	
2	2	$\delta(\boldsymbol{\eta} - 0^\circ)$	2	$\delta(\boldsymbol{\eta} - 0^\circ)$	
2	2	$\delta(\boldsymbol{\eta} - 30^\circ)$	2	$\delta(\boldsymbol{\eta} - 0^\circ)$	
2	2	$\delta(\boldsymbol{\eta} + 30^\circ)$	2	$\delta(\boldsymbol{\eta} - 0^\circ)$	

Conjecture 3: The distribution over  $\boldsymbol{\eta}$  must depend on  $\hat{\mathbf{x}}$ .



$\mathbf{x}$	$\hat{\mathbf{x}}$	$p(\boldsymbol{\eta}   \mathbf{x}, \hat{\mathbf{x}})$
2	2	$\delta(\boldsymbol{\eta} - 0^\circ)$
2	2	$\delta(\boldsymbol{\eta} - 30^\circ)$
2	2	$\delta(\boldsymbol{\eta} + 30^\circ)$
8	8	$0.5 \cdot \delta(\boldsymbol{\eta} - 0^\circ) + 0.5 \cdot \delta(\boldsymbol{\eta} - 180^\circ)$
8	8	$0.5 \cdot \delta(\boldsymbol{\eta} - 30^\circ) + 0.5 \cdot \delta(\boldsymbol{\eta} + 150^\circ)$
8	8	$0.5 \cdot \delta(\boldsymbol{\eta} + 30^\circ) + 0.5 \cdot \delta(\boldsymbol{\eta} - 150^\circ)$