



## Overview

- Widely used probabilistic models (matrix factorization, VAEs) contain *parameter symmetries* that cause approximate inference to underfit.
- We model this effect by fitting an explicitly *symmetrized* approximate posterior.
- Initial results show that this improves predictions and avoids underfitting.

## Approximate Inference -> Implicit Regularization

**Intuition:** larger models should better fit / capture structure in data.

**Observation:** in practice, variational autoencoders refuse to use extra hidden units (“component collapse”).

**Theory:** can show that variational matrix factorization (linear analogue of VAEs) ignores extra hidden units, shrinks small singular values to zero. (Nakajima et al., 2013)

This **unwanted (implicit) regularization** is caused by approximate inference – it is not present in the true Bayesian posterior!

## Symmetrized Posteriors

Classic VI: fit approximate posterior  $q$  by minimizing  $KL[q | p]$ , equivalent to maximizing an evidence lower bound (ELBO)

$$\begin{aligned} \log p(\mathbf{x}) &\geq \mathcal{L}(\theta) \\ &= \mathbb{E}_{\mathbf{z} \sim q_\theta} [\log p(\mathbf{x}, \mathbf{z})] + \mathcal{H}(q) \end{aligned}$$

Given base posterior  $q^*$ , we define the *symmetrized posterior*  $\tilde{q}$  as a uniform mixture under transformations from group  $\mathbf{G}$ :

$$\tilde{q}_\theta(\mathbf{z}) = \int_{\mathbf{T} \in \mathbf{G}} q_\theta^*(\mathbf{T}^{-1}\mathbf{z}) |\mathbf{T}^{-1}| dV(\mathbf{T})$$

Sampling interpretation: first draw  $\mathbf{z}^* \sim q^*$ , then apply a (uniformly chosen) random transformation to sample  $\mathbf{z} = \mathbf{T}\mathbf{z}^*$ .

The symmetrized posterior  $\tilde{q}$  matches symmetries of the true posterior; yields a tighter evidence bound:

$$\mathcal{L}(\tilde{q}_\theta) = \mathbb{E}_{\mathbf{z}^* \sim q_\theta^*} [\log p(\mathbf{x}, \mathbf{z}^*)] + \mathcal{H}(q^*) + KL[q^* | \tilde{q}]$$

To apply: need to compute/approximate  $KL[q^* | \tilde{q}]$  for specific symmetry group. Can do this for Gaussian  $q^*$  under orthogonal group  $O(k)$ , matching matrix factorization/VAE symmetries.

## Orthogonally symmetrized Gaussians

Symmetrizing the column space of an elementwise Gaussian matrix over the orthogonal group yields a continuous mixture of Gaussians:

$$KL[q_\theta^* | \tilde{q}_\theta] = -\mathbb{E}_{q^*} \left[ \log \int_{\mathbf{T} \in O(k)} \frac{\mathcal{N}(\mathbf{X}\mathbf{T}^T; \mathbf{M}, \Sigma)}{\mathcal{N}(\mathbf{X}; \mathbf{M}, \Sigma)} dV(\mathbf{T}) \right]$$

which decomposes as

$$\begin{aligned} &= -\mathbb{E}_{q^*} \left[ \log \int_{\mathbf{T} \in O(k)} \exp \left\{ -\frac{1}{2} \text{Tr} \left[ \mathbf{X}^T \mathbf{X} (\mathbf{T}^T \Sigma^{-1} \mathbf{T} - \Sigma^{-1}) \right] \right. \right. \\ &\quad \left. \left. + \frac{1}{2} \text{Tr} \left[ \Sigma^{-1} (\mathbf{M}^T \mathbf{X} + \mathbf{X}^T \mathbf{M}) (\mathbf{T}^T - \mathbf{I}) \right] \right\} dV(\mathbf{T}) \right] \end{aligned}$$

Taking  $\mathbf{A} = \frac{1}{2} \Sigma^{-1} (\mathbf{M}^T \mathbf{X} + \mathbf{X}^T \mathbf{M})$ , this simplifies to

$$\begin{aligned} &= \mathbb{E}_{q^*} \left[ -\log \int_{\mathbf{T} \in O(k)} \text{etr}[\mathbf{A}\mathbf{T}^T - \mathbf{A}] \right] \\ &= \mathbb{E}_{q^*} \left[ \text{Tr}[\mathbf{A}] - \log {}_0F_1 \left[ \frac{k}{2}; \frac{1}{4} \mathbf{A}\mathbf{A}^T \right] \right] \end{aligned}$$

where the hypergeometric function  ${}_0F_1$  depends only on singular values of  $\mathbf{A}$  and can be efficiently Laplace-approximated (Butler & Wood, 2003).

Intuitively, the symmetrized KL correction encourages nonzero singular values and low nullspace dimension in the mean matrix  $\mathbf{M}$ .

## Illustration: regularization from signflip symmetry

Bayesian scalar factorization:

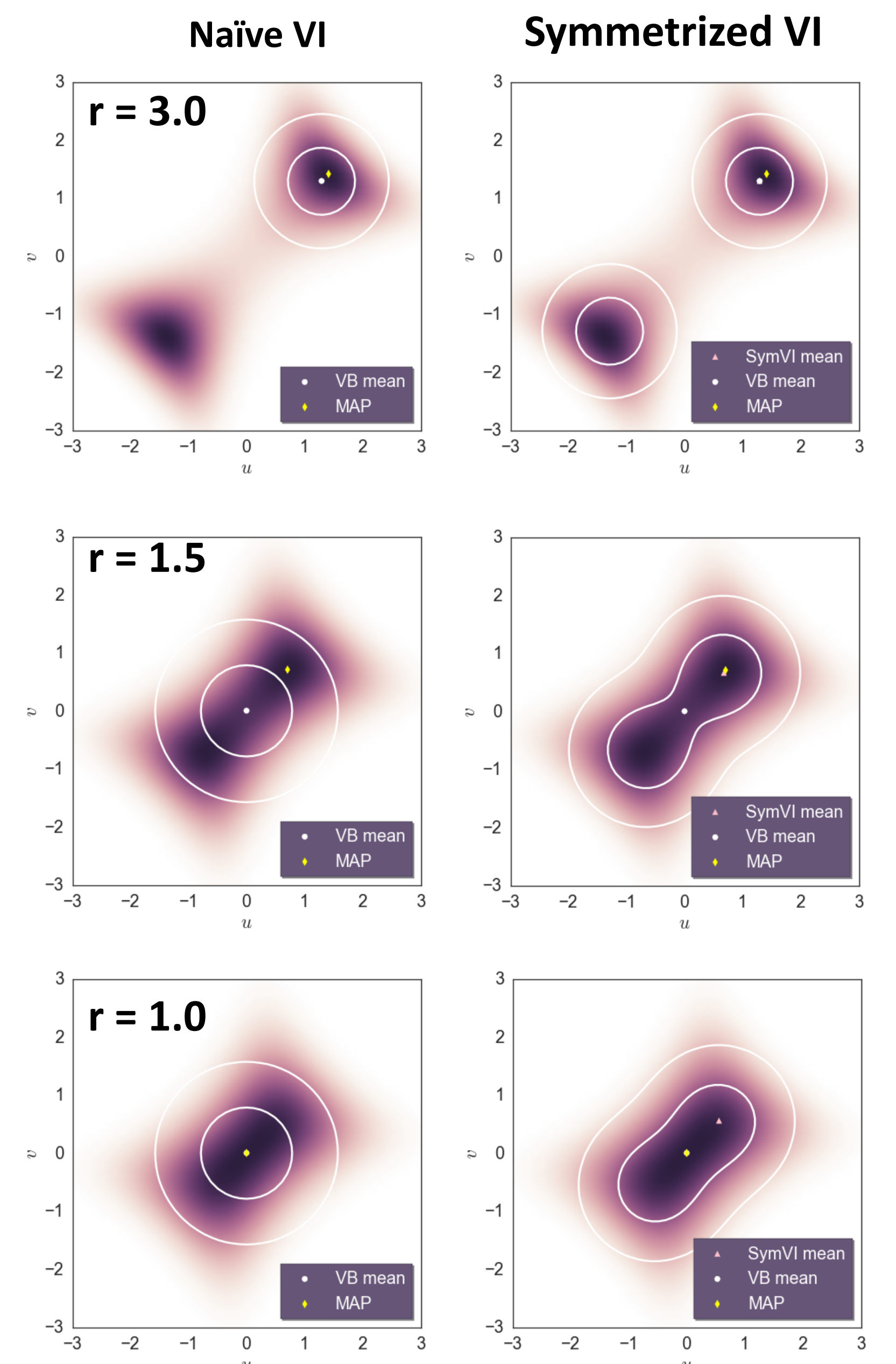
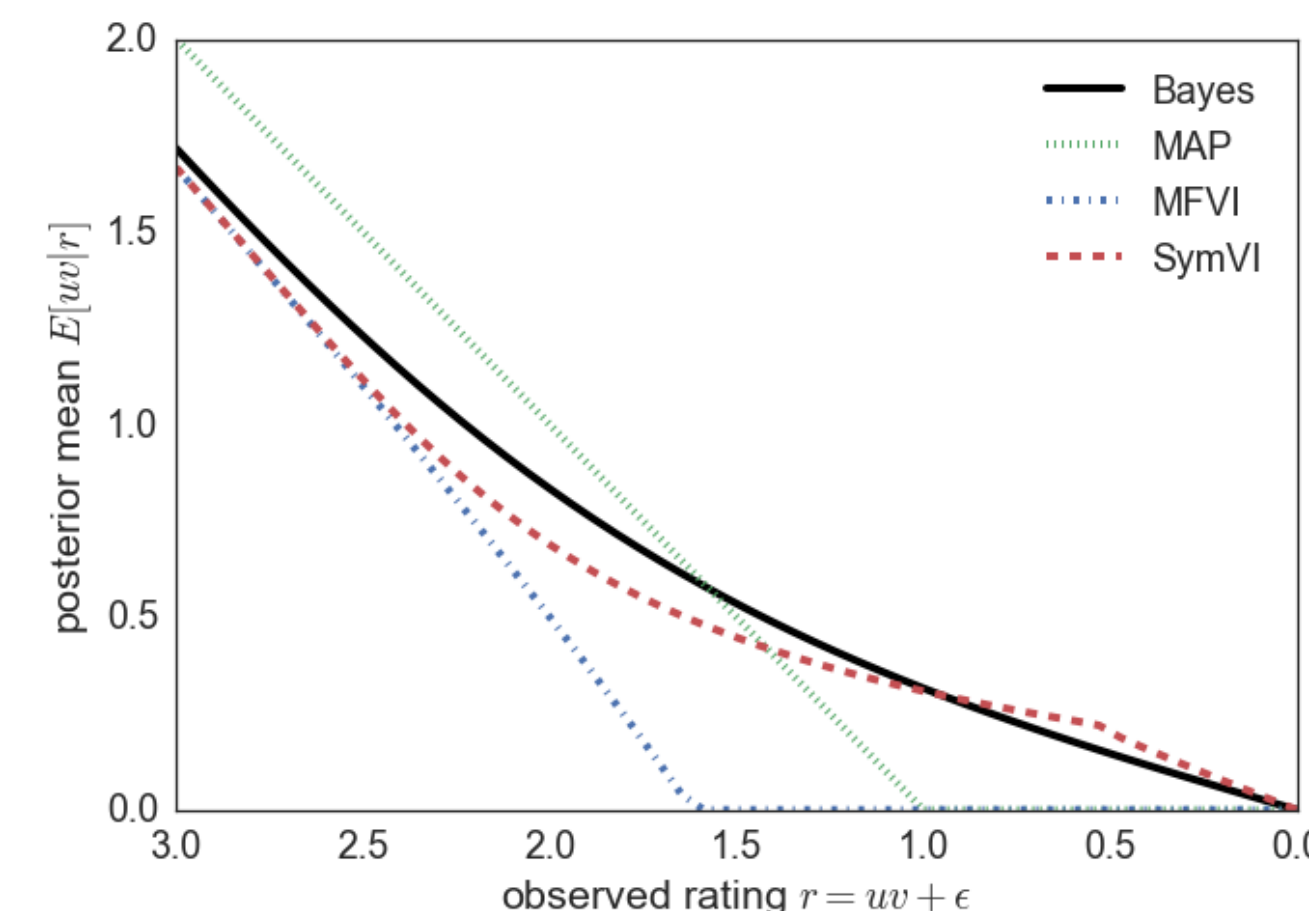
$$u, v \sim \mathcal{N}(0, 1)$$

$$\epsilon \sim \mathcal{N}(0, 1)$$

observed:  $r = uv + \epsilon$

**symmetry:**  $p(u, v | r) = p(-u, -v | r)$

**task:** predict true “rating”  $uv$



MAP and naïve VI predictions are pulled towards zero by the opposite-sign mode. Symmetrized predictions follow the true Bayes predictive mean.

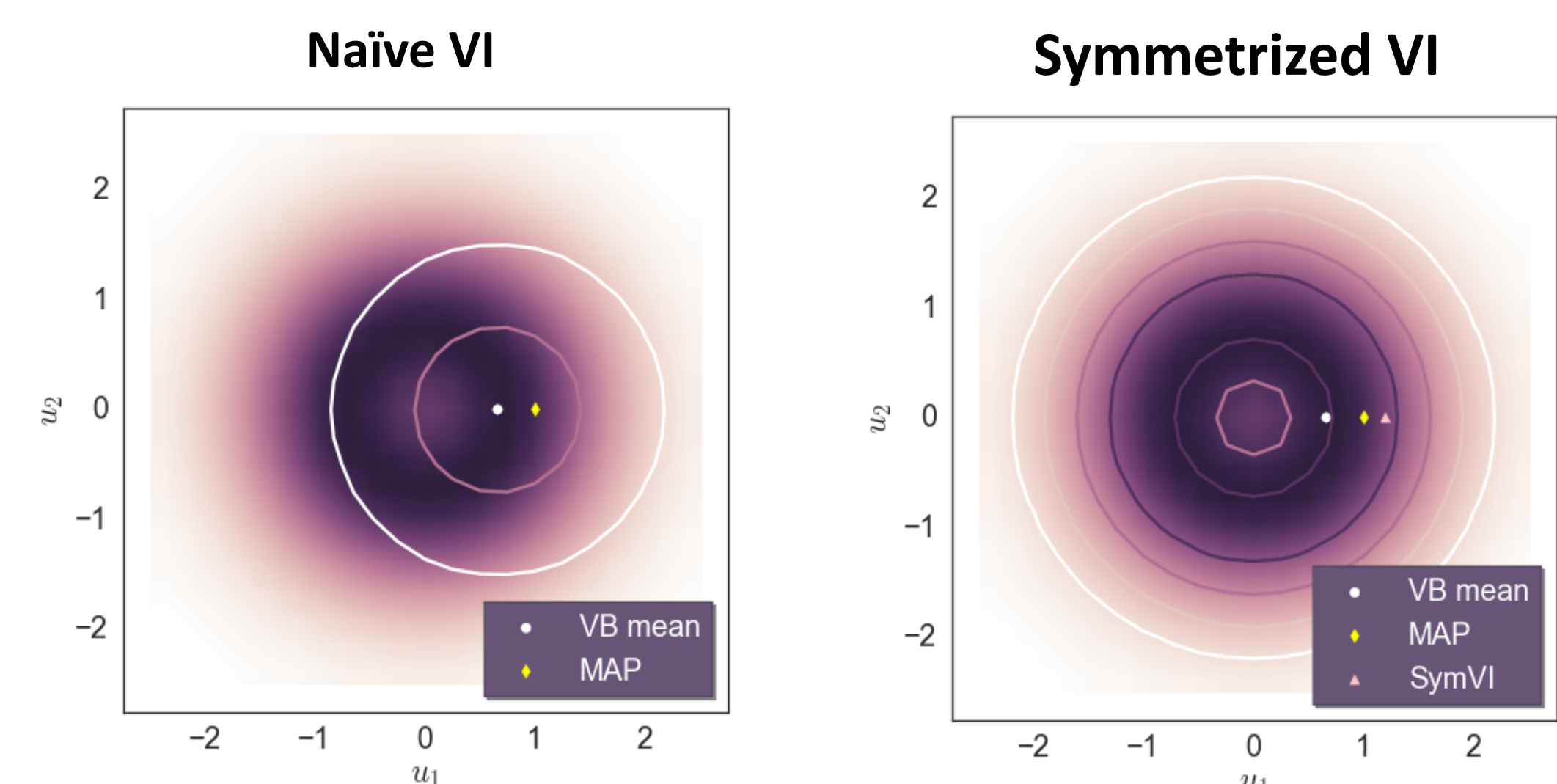
## General rotation symmetry

Bayesian matrix factorization:

$$\mathbf{R} = \mathbf{U}\mathbf{V}^T + \boldsymbol{\epsilon} = (\mathbf{U}\mathbf{T})(\mathbf{V}\mathbf{T})^T + \boldsymbol{\epsilon}$$

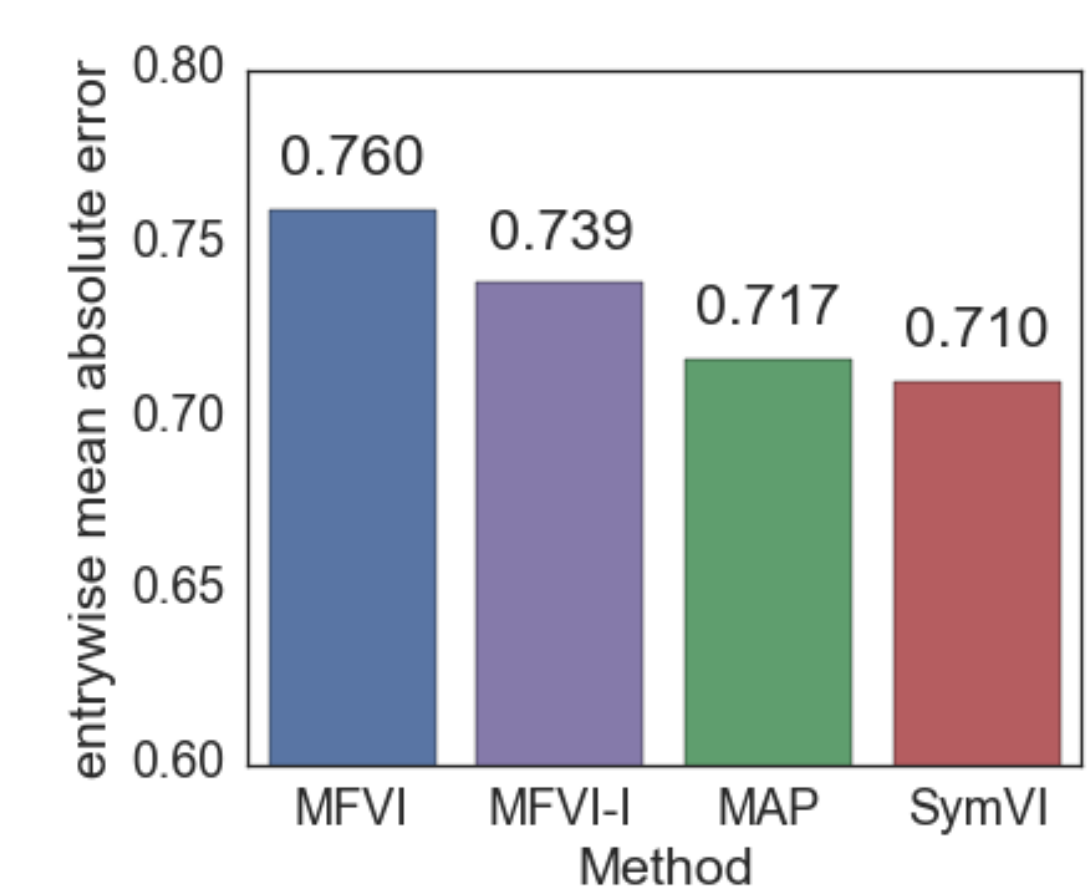
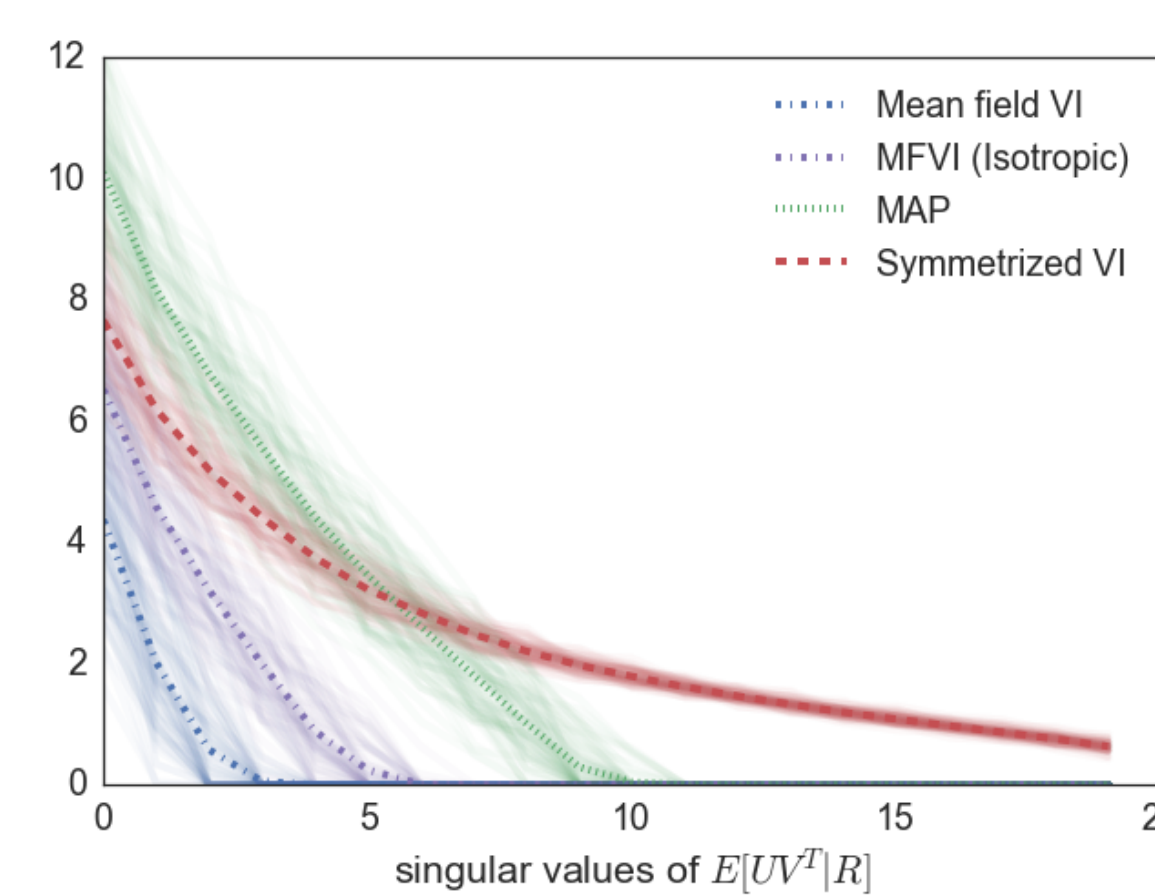
Invariant to transformation by any  $\mathbf{T}$  s.t.  $\mathbf{T}(\mathbf{T}^T) = \mathbf{I}$ , i.e., orthogonal transformations.

Can visualize in (overparameterized) case  $\mathbf{U}, \mathbf{V} \in \mathbb{R}^{1 \times 2}$ :



Naïve MAP/VI solutions are (again) shrunk towards zero. The symmetrized solution avoids shrinkage by implicitly modeling a continuous Gaussian mixture around the unit circle.

Simulations on 40 x 40 matrices with 20 latent traits:



Modeling posterior symmetries allows inference to use the full model capacity (all 20 traits)

Leads to improved predictive accuracy (recovering “true” noise-free ratings  $\mathbf{U}\mathbf{V}^T$ )

## Future/ongoing work:

- Extension to nonisotropic Gaussian  $q^*$ .
- Other expressive posterior classes (normalizing flows, autoregressive, particle-based).
- Other symmetry groups: permutation (“label switching”), translation, scaling.
- Stochastic/minibatch inference, application to VAEs.

## References

- Butler, R. W. and Wood, A. T. (2003). Laplace approximation for Bessel functions of matrix argument. *Journal of Computational and Applied Mathematics*, 155(2):359–382.
- Nakajima, S., Sugiyama, M., Babacan, S. D., and Tomioka, R. (2013). Global analytic solution of fully-observed variational Bayesian matrix factorization. *Journal of Machine Learning Research*, 14(Jan):1–37.