

SLANG: Fast Structured Covariance Approximations for Bayesian Deep Learning with Natural Gradient

Aaron Mishkin^{1,2}, Frederik Kunstner^{1,3}, Didrik Nielsen¹, Mark Schmidt², Mohammad Emtiyaz Khan¹

¹Center for Advanced Intelligence Project, RIKEN, Tokyo, Japan

²University of British Columbia, Vancouver, Canada, ³École Polytechnique Fédérale de Lausanne (EPFL), Switzerland

Introduction

Motivation:

- Uncertainty estimation is essential to make reliable decisions based on the predictions of deep models, but is computationally challenging.
- It is difficult to form even a Gaussian approximation to the posterior for large models.
- Mean-field methods reduce the computational complexity, but yield poor estimates of the uncertainty.

Contributions:

- We propose a new **stochastic, low-rank, approximate natural-gradient (SLANG)** method for Gaussian variational inference.
- Our method estimates a “low-rank plus diagonal” covariance matrix based solely on back-propagated gradients.
- SLANG is faster and more accurate than mean-field methods, and performs comparably to state-of-the-art methods.

Natural Gradient Variational Inference

Given a deep model $p(\mathcal{D}|\theta)$ with weights θ , **Gaussian Variational Inference** computes a Gaussian approximation $q(\theta) := \mathcal{N}(\theta; \mu, \Sigma)$ to the posterior by maximizing the ELBO:

$$\mathcal{L}(\mu, \Sigma) = \mathbb{E}_q[\log p(\mathcal{D}|\theta) + \log \mathcal{N}(\theta | 0, \mathbf{I}/\lambda) - \log q(\theta)],$$

Gradient-based methods optimize the ELBO using the stochastic gradient updates (t is the iteration, γ_t is the learning rate)

$$\mu_{t+1} = \mu_t - \gamma_t \hat{\nabla}_{\mu} \mathcal{L}_t, \quad \Sigma_{t+1} = \Sigma_t - \gamma_t \hat{\nabla}_{\Sigma} \mathcal{L}_t.$$

Problem: Gradient descent implicitly uses Euclidean geometry.



Natural Gradient methods do steepest descent in the space of realizable variational distributions $q(\theta)$ by optimizing on the Riemannian manifold. This gives the update

$$\mu_{t+1} = \mu_t - \beta_t \Sigma_{t+1}^{-1} \hat{\nabla}_{\mu} \mathcal{L}_t, \quad \Sigma_{t+1}^{-1} = (1 - \beta_t) \Sigma_t^{-1} + \beta_t \hat{\nabla}_{\Sigma} \mathcal{L}_t.$$

Problem: This update requires computing the Hessian.



Variational Online Gauss-Newton approximates the Hessian with the empirical Fisher Information matrix $\hat{\mathbf{G}}(\theta_t)$. This gives

$$\mu_{t+1} = \mu_t - \beta_t \Sigma_{t+1}^{-1} [\hat{\mathbf{g}}(\theta_t) + \lambda \mu_t],$$

$$\Sigma_{t+1}^{-1} = (1 - \beta_t) \Sigma_t^{-1} + \beta_t [\hat{\mathbf{G}}(\theta_t) + \lambda \mathbf{I}],$$

where $\hat{\mathbf{g}}(\theta_t)$ is the gradient and

$$\hat{\mathbf{G}}(\theta_t) = \frac{1}{M} \sum_{i=1}^M g_i(\theta_t) g_i(\theta_t)^{\top}$$

is the empirical Fisher Information matrix for $p(\mathcal{D} | \theta_t)$ computed with a minibatch of size M .

Problem: Computing and storing Σ_t is $O(D^2)$ for dense covariances.

SLANG

We approximate the covariance with a “low-rank plus diagonal” matrix

$$\Sigma_t^{-1} \approx \hat{\Sigma}_t^{-1} := \mathbf{U}_t \mathbf{U}_t^{\top} + \mathbf{D}_t,$$

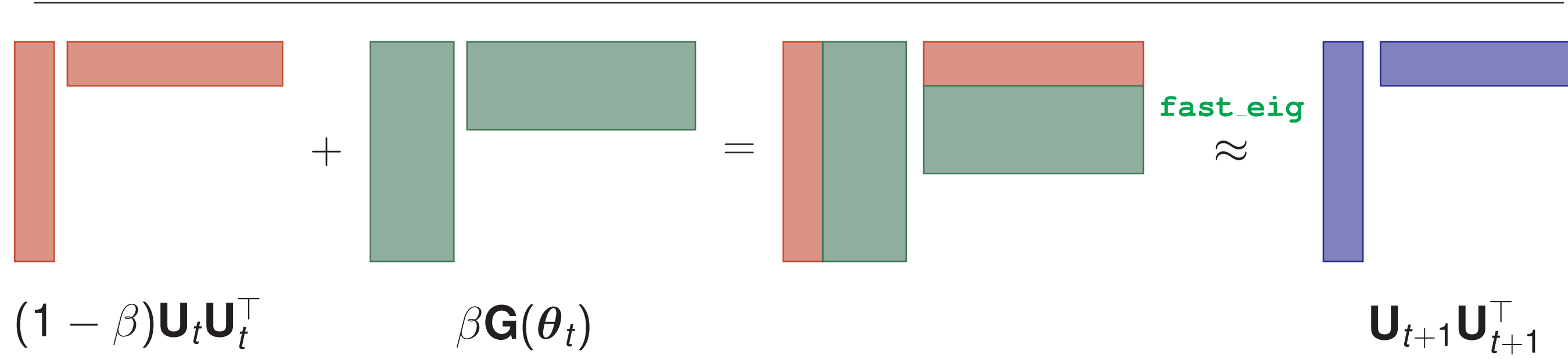
where \mathbf{U}_t is a $D \times L$ matrix and \mathbf{D}_t is diagonal. The cost of storing and inverting this matrix is linear in D which is reasonable when $L \ll D$. The approximate natural gradient update for $\hat{\Sigma}_t^{-1}$ is

$$\hat{\Sigma}_{t+1}^{-1} := \mathbf{U}_{t+1} \mathbf{U}_{t+1}^{\top} + \mathbf{D}_{t+1} \approx (1 - \beta_t) \hat{\Sigma}_t^{-1} + \beta_t [\hat{\mathbf{G}}(\theta_t) + \lambda \mathbf{I}]$$

This update may increase the rank of \mathbf{U}_{t+1} , so we project the matrix onto a L -dimensional subspace using an eigenvalue decomposition:

$$(1 - \beta_t) \hat{\Sigma}_t^{-1} + \beta_t [\hat{\mathbf{G}}(\theta_t) + \lambda \mathbf{I}] = \underbrace{(1 - \beta_t) \mathbf{U}_t \mathbf{U}_t^{\top} + \beta_t \hat{\mathbf{G}}(\theta_t)}_{\text{Rank at most } L+M} + \underbrace{(1 - \beta_t) \mathbf{D}_t + \beta_t \lambda \mathbf{I}}_{\text{Diagonal component}}$$

$$\approx \underbrace{\mathbf{Q}_{1:L} \mathbf{\Lambda}_{1:L} \mathbf{Q}_{1:L}^{\top}}_{\text{Rank } L \text{ eigendecomposition}} + \underbrace{(1 - \beta_t) \mathbf{D}_t + \beta_t \lambda \mathbf{I}}_{\text{Diagonal component}}.$$



$$(1 - \beta) \mathbf{U}_t \mathbf{U}_t^{\top} + \beta \mathbf{G}(\theta_t) = \mathbf{U}_{t+1} \mathbf{U}_{t+1}^{\top} + \mathbf{D}_{t+1}$$

The diagonal information lost in this projection is equal to

$$\Delta_D = \text{diag} \left[(1 - \beta) \mathbf{U}_t \mathbf{U}_t^{\top} + \beta_t \hat{\mathbf{G}}(\theta_t) - \mathbf{U}_{t+1} \mathbf{U}_{t+1}^{\top} \right].$$

We add this to \mathbf{D}_t as a diagonal correction. The final SLANG update is

SLANG: $\mathbf{U}_{t+1} = \mathbf{Q}_{1:L} \mathbf{\Lambda}_{1:L}^{1/2}$

$$\mathbf{D}_{t+1} = (1 - \beta) \mathbf{D}_t + \beta_t \lambda \mathbf{I} + \Delta_D.$$

$$\mu_{t+1} = \mu_t - \alpha_t \left[\mathbf{U}_{t+1} \mathbf{U}_{t+1}^{\top} + \mathbf{D}_{t+1} \right]^{-1} [\hat{\mathbf{g}}(\theta_t) + \lambda \mu_t].$$

The Algorithm

Pseudo-code for SLANG is shown in Algorithm 1. α, β are learning rates, D is denoted with a vector \mathbf{d} and \mathbf{u}_j and \mathbf{v}_j are the columns of \mathbf{U} and \mathbf{V} , respectively.

Algorithm 1: SLANG

Require: Data \mathcal{D} , hyperparameters $M, L, \lambda, \alpha, \beta$

- 1: Initialize $\mu, \mathbf{U}, \mathbf{d}$
- 2: $\delta \leftarrow (1 - \beta)$
- 3: **while** not converged **do**
- 4: $\theta \leftarrow \text{fast_sample}(\mu, \mathbf{U}, \mathbf{d})$
- 5: $\mathcal{M} \leftarrow \text{sample a minibatch}$
- 6: $[\mathbf{g}_1, \dots, \mathbf{g}_M] \leftarrow \text{backprop}(\mathcal{D}_{\mathcal{M}}, \theta)$
- 7: $\mathbf{V} \leftarrow \text{fast_eig}(\delta \mathbf{U}_1, \dots, \delta \mathbf{U}_L, \beta \mathbf{g}_1, \dots, \beta \mathbf{g}_M, L)$
- 8: $\Delta_d \leftarrow \sum_{i=1}^L \delta \mathbf{u}_i^2 + \sum_{i=1}^M \beta \mathbf{g}_i^2 - \sum_{i=1}^L \mathbf{v}_i^2$
- 9: $\mathbf{U} \leftarrow \mathbf{V}$
- 10: $\mathbf{d} \leftarrow \delta \mathbf{d} + \Delta_d + \lambda \mathbf{1}$
- 11: $\hat{\mathbf{g}} \leftarrow \sum_i \mathbf{g}_i + \lambda \mu$
- 12: $\Delta_{\mu} \leftarrow \text{fast_inverse}(\hat{\mathbf{g}}, \mathbf{U}, \mathbf{d})$
- 13: $\mu \leftarrow \mu - \alpha \Delta_{\mu}$
- 14: **end while**
- 15: **return** $\mu, \mathbf{U}, \mathbf{d}$

Algorithm 2: fast_inverse($\mathbf{g}, \mathbf{U}, \mathbf{d}$)

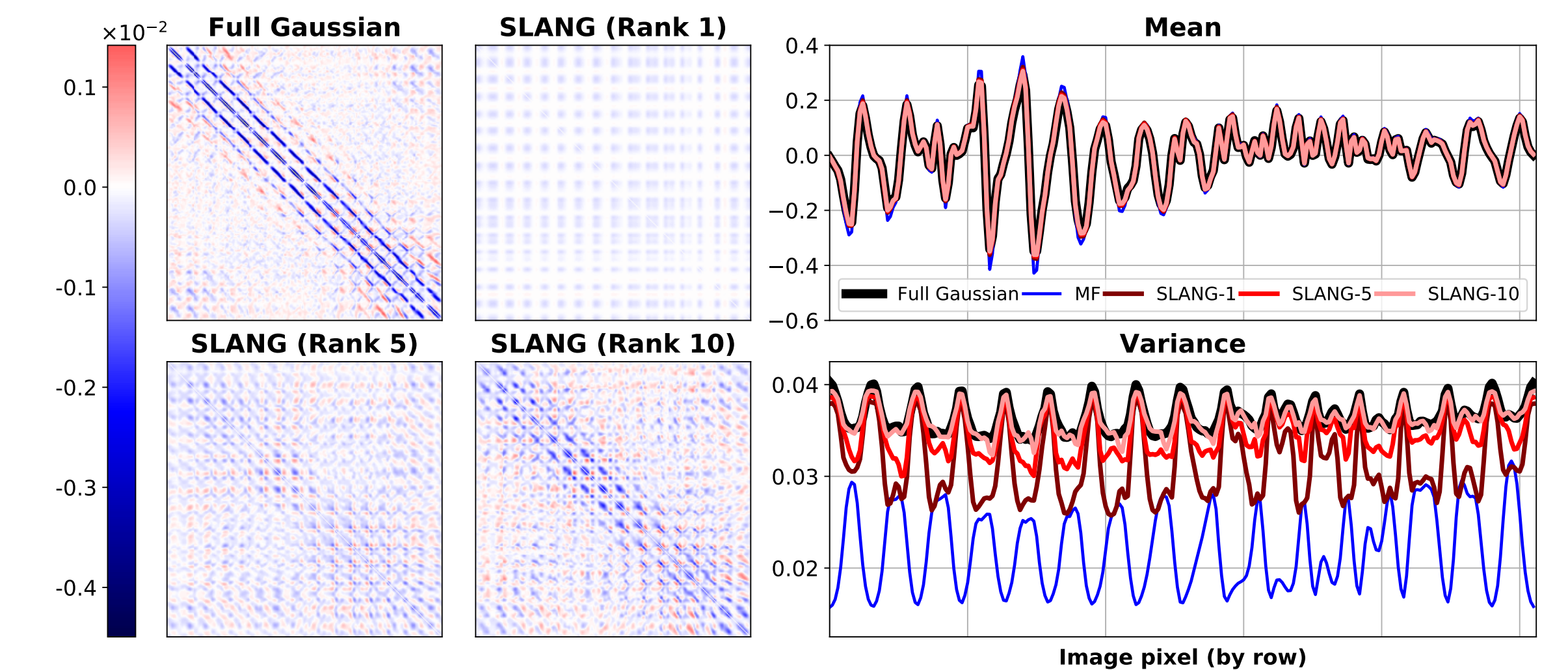
- 1: $\mathbf{A} \leftarrow (\mathbf{I}_L + \mathbf{U}^{\top} \mathbf{d}^{-1} \mathbf{U})^{-1}$
- 2: $\mathbf{y} \leftarrow \mathbf{d}^{-1} \mathbf{g} - \mathbf{d}^{-1} \mathbf{U} \mathbf{A} \mathbf{U}^{\top} \mathbf{d}^{-1} \mathbf{g}$
- 3: **return** \mathbf{y}

Algorithm 3: fast_sample($\mu, \mathbf{U}, \mathbf{d}$)

- 1: $\epsilon \sim \mathcal{N}(0, \mathbf{I}_D)$
- 2: $\mathbf{V} \leftarrow \mathbf{d}^{-1/2} \odot \mathbf{U}$
- 3: $\mathbf{A} \leftarrow \text{Cholesky}(\mathbf{V}^{\top} \mathbf{V})$
- 4: $\mathbf{B} \leftarrow \text{Cholesky}(\mathbf{I}_L + \mathbf{V}^{\top} \mathbf{V})$
- 5: $\mathbf{C} \leftarrow \mathbf{A}^{-\top} (\mathbf{B} - \mathbf{I}_L) \mathbf{A}^{-1}$
- 6: $\mathbf{K} \leftarrow (\mathbf{C} + \mathbf{V}^{\top} \mathbf{V})^{-1}$
- 7: $\mathbf{y} \leftarrow \mathbf{d}^{-1/2} \epsilon - \mathbf{V} \mathbf{K} \mathbf{V}^{\top} \epsilon$
- 8: **return** $\mu + \mathbf{y}$

Results

Covariance Structure for Logistic Regression on USPS



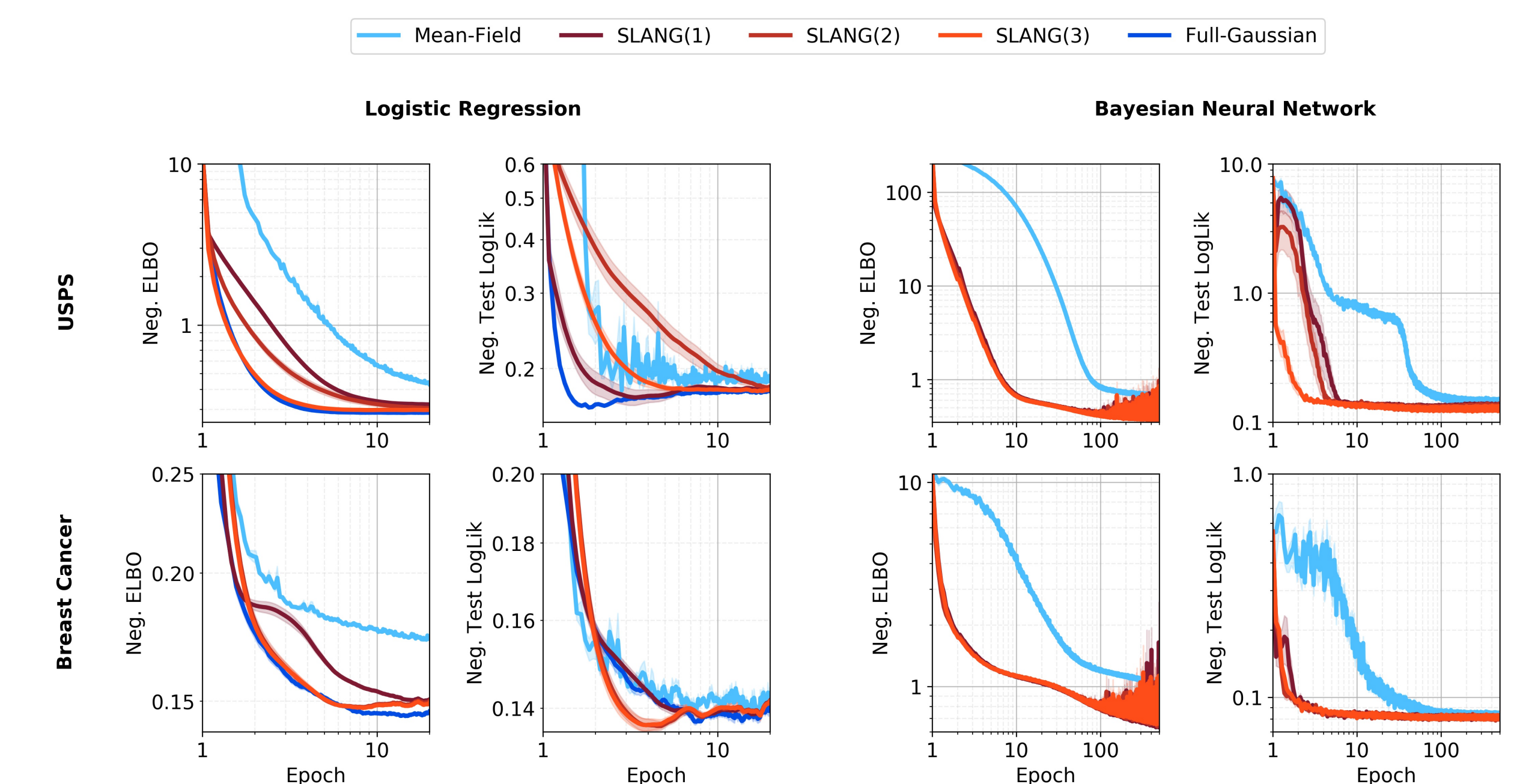
- SLANG doesn't underestimate variance like mean-field methods.

Logistic Regression Results

Dataset	Metrics	Mean-Field Methods			SLANG			Full Gaussian		
		EF	Hess.	Exact	L = 1	L = 5	L = 10	EF	Hess.	Exact
Australian	NLL	0.348	0.347	0.341	0.342	0.339	0.338	0.340	0.339	0.338
	KL ($\times 10^4$)	2.240	2.030	0.195	0.033	0.008	0.002	0.000	0.000	0.000
a1a	NLL	0.339	0.339	0.339	0.339	0.339	0.339	0.339	0.339	0.339
	KL ($\times 10^2$)	2.590	2.208	1.295	0.305	0.173	0.118	0.014	0.000	0.000
USPS	NLL	0.139	0.139	0.138	0.132	0.132	0.131	0.131	0.130	0.130
3vs5	KL ($\times 10^1$)	7.684	7.188	7.083	1.492	0.755	0.448	0.180	0.001	0.000

- SLANG performs similarly to full-Gaussian methods at test time.

Convergence Experiments



- SLANG converges faster than mean-field methods for logistic regression and BNNs.

Bayesian Neural Networks Results:

Dataset	BBB	Test RMSE			Test log-likelihood		
		Dropout	SLANG		BBB	Dropout	SLANG
Boston	3.43 \pm 0.20	2.97 \pm 0.19	3.21 \pm 0.19		-2.66 \pm 0.06	-2.46 \pm 0.06	-2.58 \pm 0.05
Concrete	6.16 \pm 0.13	5.23 \pm 0.12	5.58 \pm 0.19		-3.25 \pm 0.02	-3.04 \pm 0.02	-3.13 \pm 0.03
Energy	0.97 \pm 0.09	1.66 \pm 0.04	0.64 \pm 0.03		-1.45 \pm 0.10	-1.99 \pm 0.02	-1.12 \pm 0.01
Kin8nm	0.08 \pm 0.00	0.10 \pm 0.00	0.08 \pm 0.00		1.07 \pm 0.00	0.95 \pm 0.01	1.06 \pm 0.00
Naval	0.00 \pm 0.00	0.01 \pm 0.00	0.00 \pm 0.00		4.61 \pm 0.01	3.80 \pm 0.01	4.76 \pm 0.00
Power	4.21 \pm 0.03	4.02 \pm 0.04	4.16 \pm 0.04		-2.86 \pm 0.01	-2.80 \pm 0.01	-2.84 \pm 0.01

- Performance on BNNs is comparable to Bayesian Dropout and Bayes-by-Backprop.