UNIVERSITY OF COPENHAGEN DEPARTMENT OF COMPUTER SCIENCE



Prediction

y_pred = einstein.predict(state, X_test)['y'] print(f"Accuracy: {np.mean(y_pred == y_test)*100:.2f} %")

callbacks=[Progbar()])

Integrated Stein VI Theory

Bayesian Variational Inference

- Goal: Given prior p(z) over model parameters z and likelihood p(x|z)over data x infer posterior: $p(z|x) = Z^{-1}p(x|z)p(z)$
- Problem: Normalization constant $Z = \int p(x|z) dx$ is intractable.
- Solution: Variational Inference
- Posit easy to evaluate approximate distribution q(z)
- Minimize the Kullback-Leibler (KL) divergence $D_{KL}(q \parallel p)$ which can be done without explicitly calculating Z.
- Stein VI: VI (Liu and Wang 2016) Use particles $\{z_i\}_{i=1}^N$ for approximation, so $q(z) = \sum_{i} \delta_{z_{i}}(z)$ and minimize KL divergence by iterating Stein forces S(z).

Two Forces of EinStein VI

- Assumption: We have a negative loss function $\mathcal{L}(\boldsymbol{\theta})$ we would like to maximize.
- **Example:** Evidence Lower Bound (ELBO; Kingma and Welling, 2014).
- **Desiderata:** We would like to have flexibility to capture multi-modality for parameters $\boldsymbol{\theta}$ and quantify uncertainty on it.
- Idea: Use ELBO-within-Stein inference (Nalisnick, 2017) which can be factored into an attractive force: $S^{+}(\boldsymbol{\theta}) = \mathbb{E}_{\boldsymbol{\vartheta} \sim q}[k(\boldsymbol{\theta}, \boldsymbol{\vartheta}) \nabla_{\boldsymbol{\theta}} L(\boldsymbol{\theta})]$ and a repulsive force:
 - $S^{-}(\boldsymbol{\theta}) = \mathbb{E}_{\boldsymbol{\vartheta} \sim q} [\nabla_{\boldsymbol{\theta}} k(\boldsymbol{\theta}, \boldsymbol{\vartheta})]$ where $k: \mathbb{R}^d \times \mathbb{R}^d \to \mathbb{R}$ is a statistical kernel like RBF.

References

D. Phan, N. Pradhan, and M. Jankowiak. Composable effects for flexible and accelerated probabilistic programming in NumPyro. Program Transformational inference algorithm. NeurIPS, 2018. | E. Nalisnick. Variational inference algorithm. NeurIPS, 2018. | D. Wang. Stein variational inference with Stein mixtures. In Advances in Approximate Bayesian Inference, 2017. | D. Wang. and Q. Liu. Nonlinear stein variational gradient descent for learning diversified mixture models. ICML 2019. | D. Wang, Z. Tang, C. Bajaj, and Q. Liu. Stein variational gradient descent with matrix-valued kernels. NeurIPS 2019, 2019. | M. D. Hoffman, P. Sountsov, J. Bajaj, and Q. Liu. Stein variational gradient descent with matrix-valued kernels. NeurIPS 2019, 2019. | M. D. Hoffman, P. Sountsov, J. Bajaj, and Q. Liu. Stein variational gradient descent with matrix-valued kernels. NeurIPS 2019, 2019. | M. D. Hoffman, P. Sountsov, J. Bajaj, and Q. Liu. Stein variational gradient descent with matrix-valued kernels. NeurIPS 2019, 2019. | M. D. Hoffman, P. Sountsov, J. Bajaj, and Q. Liu. Stein variational gradient descent with matrix-valued kernels. NeurIPS 2019, 2019. | M. D. Hoffman, P. Sountsov, J. Bajaj, and Q. Liu. Stein variational gradient descent with matrix-valued kernels. NeurIPS 2019, 2019. | M. D. Hoffman, P. Sountsov, J. Bajaj, and Q. Liu. Stein variational gradient descent with matrix-valued kernels. NeurIPS 2019, 2019. | M. D. Hoffman, P. Sountsov, J. Bajaj, and Q. Liu. Stein variational gradient descent with matrix-valued kernels. NeurIPS 2019, 2019. | M. D. Hoffman, P. Sountsov, J. Bajaj, and Q. Liu. Stein variational gradient descent with matrix-valued kernels. NeurIPS 2019, 2019. | M. D. Hoffman, P. Sountsov, J. Bajaj, and Q. Liu. Stein variational gradient descent with matrix-valued kernels. NeurIPS 2019, 2019. | M. D. Hoffman, P. Sountsov, J. Bajaj, and Q. Liu. Stein variational gradient descent with matrix-valued kernels. NeurIPS 2019, 2019. | M. D. Hoffman, P. Sountsov, J. Bajaj, and Q. Liu. Stein variational gradient descent with matrix-valued kernels. NeurIPS 2019, 2019. | M. D. Hoffman, P. Sountsov, J. Bajaj, and Q. Liu. Stein variational gradient descent with matrix-valued kernels. NeurIPS 2019, 2019. | M. D. Hoffman, P. Sountsov, J. Bajaj, and Q. Liu. Stein variational gradient descent with matrix-valued kernels. NeurIPS 2019, 2019. | M. D. Hoffman, P. Sountsov, J. Bajaj, and Q. Liu. Stein Dillon, I. Langmore, D. Tran, and S. Vasudevan. Neutra-lizing bad geometry in hamiltonian monte carlo using neural transport. 2017. M. Jankowiak and T. Karaletsos. Pathwise derivatives for multivariate distributions. AISTATS 2019, 2019.

EinStein VI: General and Integrated Stein Variational Inference in NumPyro

Ahmad Salim Al-Sibahi (UCPH), Ola Rønning (UCPH), Christophe Ley (UGent), Thomas Hamelryck (UCPH) Implementation Challenges

Einstein VI is a NumPyro (Phan 2019) library that integrates the latest developments of Stein VI (Liu and Wang 2016):

- ✓ Compositional library supporting various loss functions (ELBO, Rényi ELBO, custom loss, etc.), automatic marginalization of discrete variables (Obermeyer et al 2019), and deep learning.
- ✓ Inference based on <u>ELBO</u>-with<u>in</u>-<u>Stein</u> core (Nalisnick 2017) with support for non-linear optimization (Wang and Liu 2019), a wealth of kernels that include matrix-valued ones (Wang et al. 2019) for graphical models and higher-order optimization.
- ✓ Learnable transforms on the parameter space, including triangular transforms (Parno and Marzouk 2018) and neural transforms (Hoffman et al. 2018)

Matrix-Valued Kernels

- Idea: (Wang et al. 2019) Replace scalar-valued kernel k with matrixvalued kernel $K : \mathbb{R}^d \times \mathbb{R}^d \to \mathbb{R}^{d \times d}$
- **Update:** Change the attractive force: $S^{+}(\boldsymbol{\theta}) = \mathbb{E}_{\boldsymbol{\vartheta} \sim q}[K(\boldsymbol{\theta}, \boldsymbol{\vartheta}) \nabla_{\boldsymbol{\theta}} \mathcal{L}(\boldsymbol{\theta})]$ and repulsive force:
- $S^{-}(z) = \mathbb{E}_{z' \sim q(z)}[K(z, z')\nabla_{z}]$ • Advantages: Allow pre-conditioning using second-order matrices like the Hessian or Fisher information. Allow factorization of *K* into local kernels $\{K^{(\ell)}\}_{\rho}$ based on independencies given by a graphical model.

Non-linear Stein

• Idea: (Wang and Liu 2019) Allow scaling of repulsive force $S^{-}(\theta)$ by constant

Algorithm

Input: Classical parameters ϕ and ϕ' , Stein parameters $\{\theta_i\}_i$, model $p_{\phi}(\mathbf{z}, \mathbf{x})_i$ guide $q_{\theta,\phi'}(\mathbf{z})$, loss \mathcal{L} , kernel interface KI. **Output:** Parameter changes based on classical VI ($\Delta \phi$, $\Delta \phi'$) and Stein VI forces $(\{\Delta \boldsymbol{\theta}_i\}_i)$. $\Delta \boldsymbol{\phi} \leftarrow \mathbb{E}_{\boldsymbol{\theta}} [\nabla_{\boldsymbol{\phi}} \mathcal{L}(p_{\boldsymbol{\phi}}, q_{\boldsymbol{\theta}, \boldsymbol{\phi}'})]$ $\Delta \boldsymbol{\phi}' \leftarrow \mathbb{E}_{\boldsymbol{\theta}} [\nabla_{\boldsymbol{\phi}'} \mathcal{L}(p_{\boldsymbol{\phi}}, q_{\boldsymbol{\theta}, \boldsymbol{\phi}'})]$ $\{\mathbf{a}_i\}_i \leftarrow \text{PYTREEFLATTEN}(\{\boldsymbol{\theta}_i\}_i)$ $k \leftarrow \mathrm{KI}(\{\mathbf{a}_i\}_i)$ $\{\Delta \mathbf{a}_i\}_i \leftarrow \mathrm{VMAP}(\{\mathbf{a}_i\}_i, \mathbf{a}_i \mapsto \sum_{\mathbf{a}_i} k(\mathbf{a}_j, \mathbf{a}_i) \nabla_{\mathbf{a}_i} \mathcal{L}(p_{\phi}, q_{\mathrm{PYTREERESTORE}(\mathbf{a}), \phi'}) +$ $\nabla_{\mathbf{a}_i} k(\mathbf{a}_j, \mathbf{a}_i))$ $\{\Delta \boldsymbol{\theta}_i\}_i \leftarrow \text{PYTREERESTORE}(\{\Delta \mathbf{a}_i\}_i)$ return $\Delta \phi$, $\Delta \phi'$, $\{\Delta \theta_i\}_i$

Parameter Transforms

Assumption: We can write parameters for loss \mathcal{L}_{θ} as $\theta = \mathcal{T}(\phi)$. Idea: We can reparametrize the Stein force:

NumPyro code example on the right

Experiments and Results



Neal's Funnel

 $S_{\Phi}(\phi) = S_{\Theta}(\mathcal{T}(\phi))(\nabla_{\phi}\mathcal{T}(\phi))^{\mathsf{T}}$

 $y \sim \mathcal{N}(0,3) \quad x \sim \mathcal{N}(0, \exp(\frac{y}{2}))$ **Problem:** As y becomes negative, x becomes harder to sample using MCMC. **Result:** EinStein VI can be seen to perform well in the above figure!

Double Moon Model

• **Problem:** Multi-modal distribution which can be hard to fit using traditional VI

def guide(n_flows=3, h_dim=[2,2]):

 $flows = flows(n_flows, 2, h_dim)$

particle = param('p', array([0.,0.]),

numpyro.sample('x', Delta(particle))

ComposeTrans(flows))

• **Result**: EinStein VI with neural transform (3 autoregressive flows of hidden dimensions (2,2)) works well!



Stein Mixture Latent Dirichlet Allocation (SM-LDA)

- following process:
- $\boldsymbol{\theta} \sim \text{Dir}(\alpha) \quad z_n \sim^{\text{iid}} \text{Cat}(\boldsymbol{\theta}) \quad w_n \sim^{\text{iid}} \text{Cat}(\varphi_{z_n}) \quad n \in \{1...N\}$ • **Problem:** Discrete latent variable Z_n makes model nondifferentiable
- Implementation: MLP (100 hidden dimensions), 20 topics, 5 Stein particles
- **Result:** EinStein VI works well with automatic marginalization provided by NumPyro (Obermeyer et al. 2019)

Negative Krishnar

6.85

existing results!

Implementing EinStein VI required dealing with the following challenges:

- **Re-initialiable guides** We need to re-run the initialization procedure for each Stein particle, so that they get different values for inference. **WrappedGuide** in the intro code example allows this support.
- NumPyro and JAX integration We rely on NumPyro for writing the probabilistic programs and integrated with their
- interface. The algorithm (left) needs JAX functionality such as vmap to calculate particle loss in parallel without changing the user code and PyTrees to map back
- and forth between parameters encoded as Python structures and their monolithic vectorized particle form that can capture correlation between particles.
- Kernels We needed a modular interface to support new kernels, based on particle information and log-probability.

• **Goal:** Infer topics $\boldsymbol{\theta}$ from document represented as bags of words **w**. Each document is generated according to the

Stein Mixture Deep Markov Model (SM-DMM)

est log-likelihood for JSB Chorales dataset using SM-DMM		
nd Shalit 2017	Jankowiak and Karaletsos 2019	EinStein VI (<i>Ours)</i>
	6.82	6.67

Goal: Sequential model inspired by HMMs but with deep neural networks for transitions and emissions. GRU guide. **Problem:** High-dimensional guide, over 500.000 parameters **Implementation:** 5 Stein particles, Adam optimizer (lr: 1e-5) **Result:** EinStein VI scales well and performs better than