

Scalable Training of Inference Networks for Gaussian-Process Models

Jiaxin Shi

Tsinghua University

Joint work with

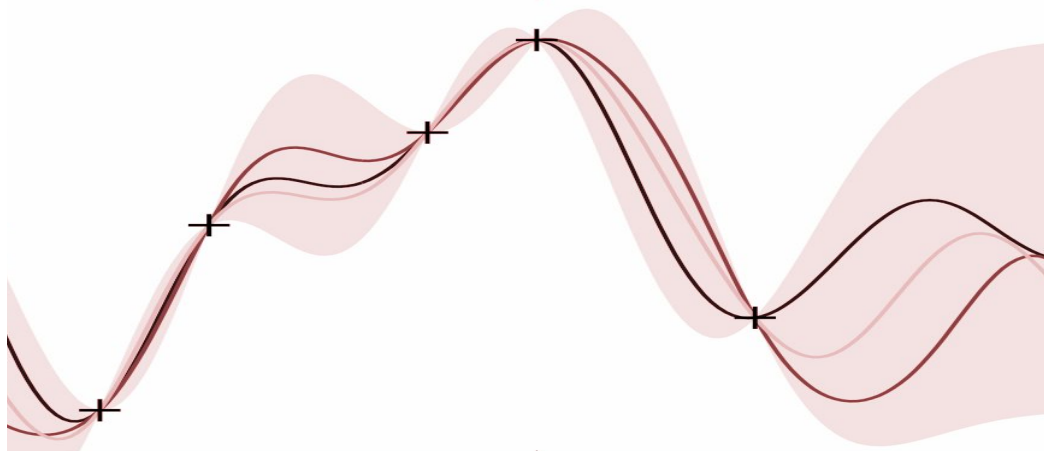


Mohammad Emtiyaz Khan



Jun Zhu

Gaussian Process

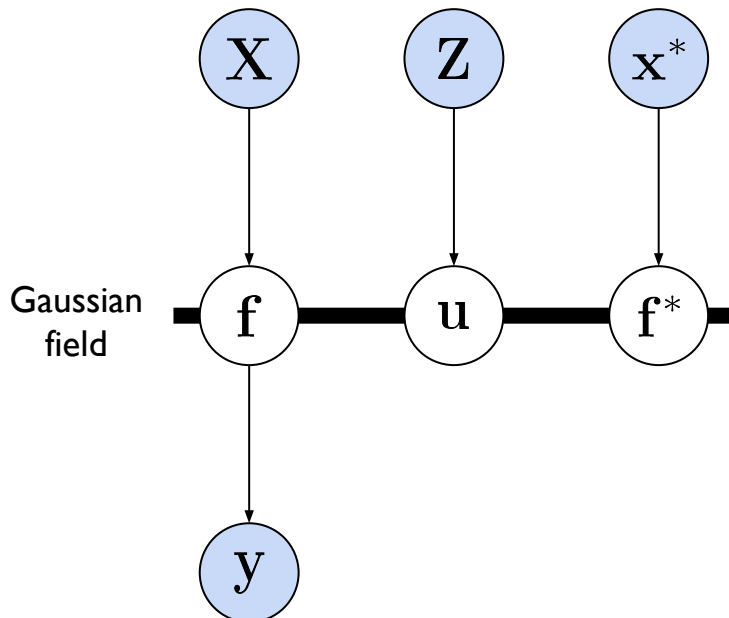


$$f(\mathbf{x}) \sim \mathcal{GP}(m(\mathbf{x}), k(\mathbf{x}, \mathbf{x}'))$$

mean function

covariance function / kernel

inducing points



Posterior inference

$$p(\mathbf{f}, \mathbf{f}^* | \mathbf{y}) \propto p(\mathbf{f}, \mathbf{f}^*) p(\mathbf{y} | \mathbf{f})$$

$\mathcal{O}(N^3)$ complexity, conjugate likelihoods

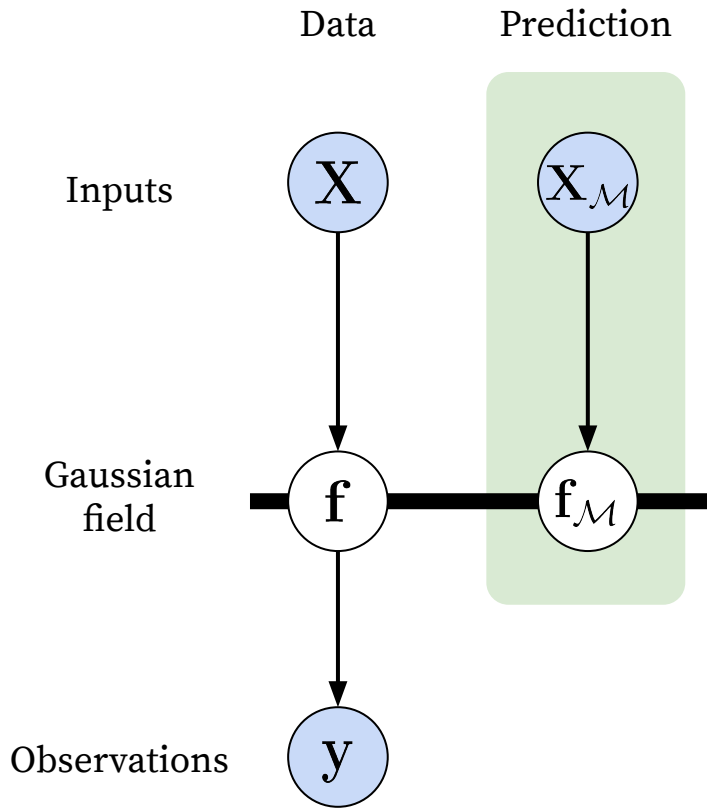
Sparse variational GP [Titsias, 09; Hensman et al., 13]

$$q(\mathbf{f}, \mathbf{f}^*, \mathbf{u}) := q(\mathbf{u}) p(\mathbf{f}, \mathbf{f}^* | \mathbf{u})$$

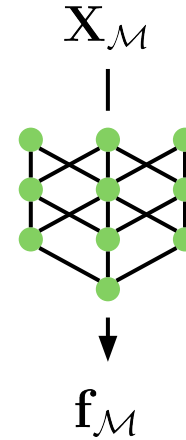
$$\mathcal{L}(q, \mathbf{Z}) := \mathbb{E}_{q(\mathbf{u}) p(\mathbf{f} | \mathbf{u})} [\log p(\mathbf{y} | \mathbf{f})] - \text{KL}[q(\mathbf{u}) || p(\mathbf{u})]$$

Inference Networks for GP Models

Remove sparse assumption



\approx

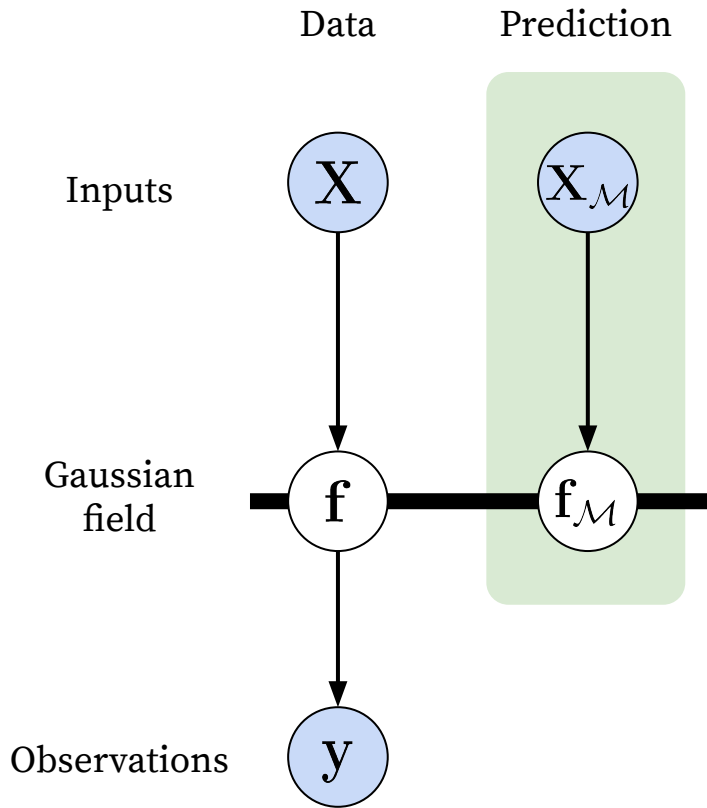


Inference network

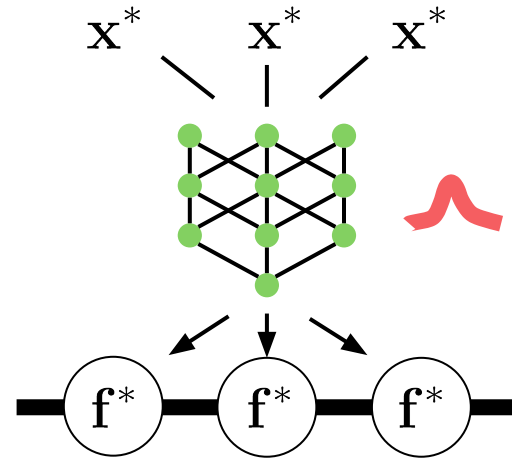
$$q(\mathbf{f}_{\mathcal{M}}) = \mathcal{N}(\mathbf{f}_{\mathcal{M}} | \mu_{\mathcal{M}}, \Sigma_{\mathcal{M}})$$

Inference Networks for GP Models

Remove sparse assumption



\approx



Inference network

$$q(\mathbf{f}_{\mathcal{M}}) = \mathcal{N}(\mathbf{f}_{\mathcal{M}} | \mu_{\mathcal{M}}, \Sigma_{\mathcal{M}})$$

Examples of Inference Networks

- Bayesian neural networks: [Sun et al., 19]
 - **intractable** output density

function space $f(\mathbf{x}) \sim \mathcal{GP}(0, k(\mathbf{x}, \mathbf{x}'))$

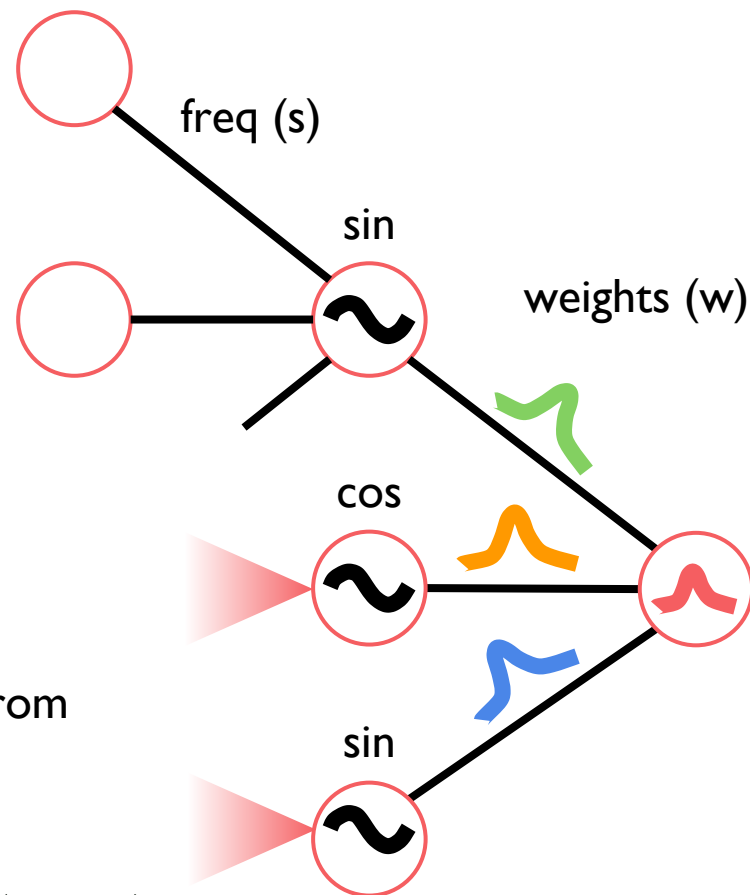


weight space $f(\mathbf{x}) = \mathbf{w}^\top \phi(\mathbf{x}), \quad \mathbf{w} \sim \mathcal{N}(\mathbf{0}, \mathbf{I})$

- Inference network architecture can be derived from the weight-space posterior

$q(f) : f(\mathbf{x}) = \mathbf{w}^\top \phi(\mathbf{x}) + \xi_\theta(\mathbf{x}), \quad \mathbf{w} \sim \mathcal{N}(\mathbf{m}, \mathbf{V})$

- Random feature expansions [Cutajar, et al., 18]
- Deep neural nets



Minibatch Training is Difficult

Functional Variational Bayesian Neural Networks (Sun et al., 19)

Measurement points

- Consider matching variational and true posterior processes at arbitrary $\mathbf{x}^{\mathcal{M}}$

$$\text{KL}[q(\mathbf{f}^{\mathcal{M}}) \| p(\mathbf{f}^{\mathcal{M}} | \mathbf{y})] \leq \text{KL}[q_{\phi}(\mathbf{f}^{\mathcal{M}}, \mathbf{f}) \| p(\mathbf{f}^{\mathcal{M}}, \mathbf{f} | \mathbf{y})]$$

- Full batch fELBO

$$\mathcal{L}_{\mathbf{x}^{\mathcal{M}}, \mathbf{x}}(q) = \log p(\mathcal{D}) - \text{KL}[q(\mathbf{f}^{\mathcal{M}}, \mathbf{f}) \| p(\mathbf{f}^{\mathcal{M}}, \mathbf{f} | \mathbf{y})].$$

$$= \sum_{(\mathbf{x}, y) \in \mathcal{D}} \mathbb{E}_{q_{\phi}} [\log p(y | f(\mathbf{x}))] - \text{KL}[q(\mathbf{f}^{\mathcal{M}}, \mathbf{f}) \| p(\mathbf{f}^{\mathcal{M}}, \mathbf{f})]$$

- Practical fELBO

$$\frac{1}{|\mathcal{D}_s|} \sum_{(\mathbf{x}, y) \in \mathcal{D}_s} \mathbb{E}_{q_{\phi}} [\log p(y | f(\mathbf{x}))] - \lambda \text{KL}[q(\mathbf{f}^{\mathcal{D}_s}, \mathbf{f}^{\mathcal{M}}) \| p(\mathbf{f}^{\mathcal{D}_s}, \mathbf{f}^{\mathcal{M}})].$$

- This objective is doing **improper** minibatch for the KL divergence term

Scalable Training of Inference Networks for GP Models

Stochastic, functional mirror descent

- work with the **functional density** directly [Dai et al., 16; Cheng & Boots, 16]
 - natural gradient in the density space
 - minibatch approximation with **stochastic** functional gradient

$$q_{t+1} = \operatorname{argmax}_q \int \hat{\partial} \mathcal{L}(q_t) q(f) df - \frac{1}{\beta_t} \text{KL} [q \| q_t]$$

- closed-form solution as an **adaptive Bayesian filter**

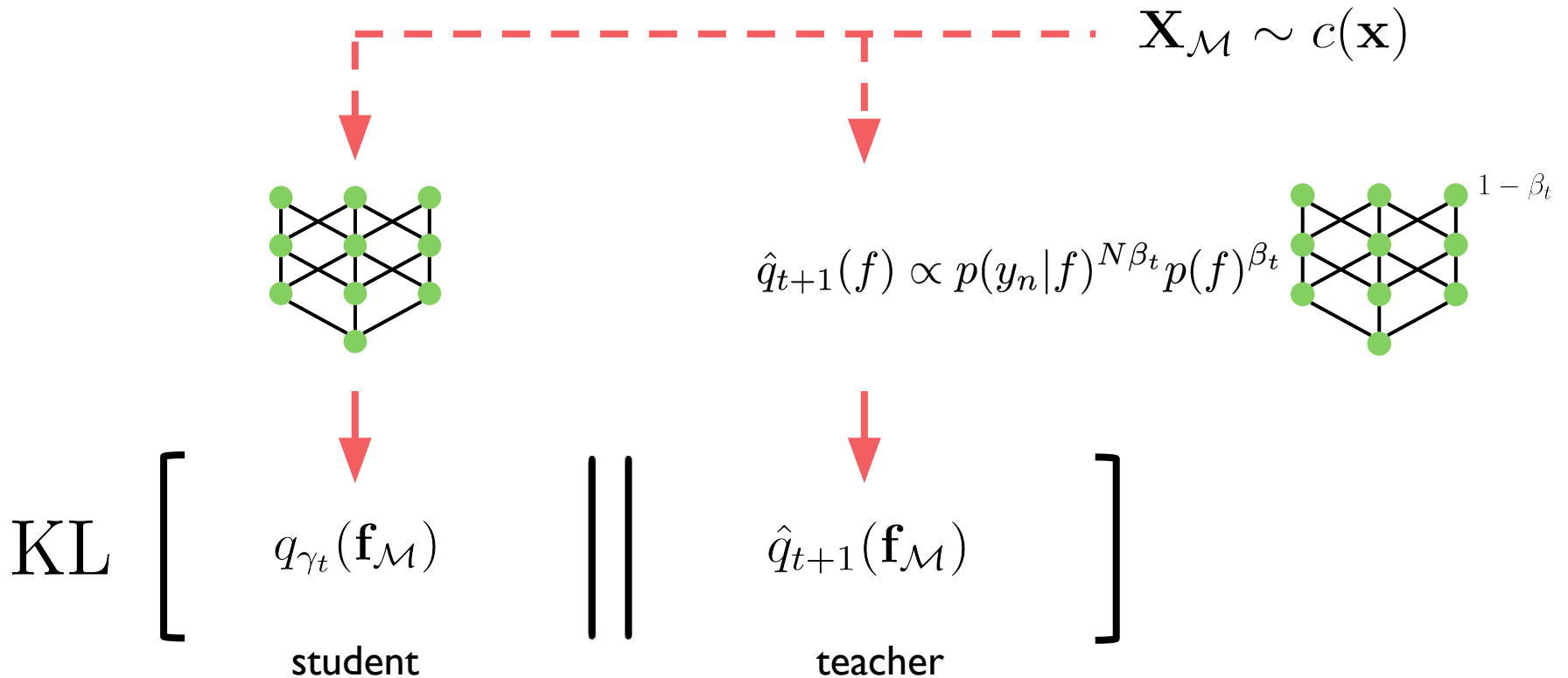
$$q_{t+1}(f) \propto \underbrace{p(y_n | f)^{N\beta_t}}_{\text{seeing next data point}} \underbrace{p(f)^{\beta_t} q_t(f)^{1-\beta_t}}_{\text{adapted prior}}$$

- sequentially applying Bayes' rule is the most **natural** gradient
 - in conjugate models: equivalent to natural gradient for exponential families

[Raskutti & Mukherjee, 13; Khan & Lin, 17]

Scalable Training of Inference Networks for GP Models

Minibatch training of inference networks



- an idea from filtering: **bootstrap**
 - similar idea: **temporal difference** (TD) learning with function approximation

Scalable Training of Inference Networks for GP Models

Minibatch training of inference networks

- **(Gaussian likelihood case)** closed-form marginals of $\hat{q}_{t+1}(f)$ at locations $\mathbf{X}_{\mathcal{M}}$
 - equivalent to **GP regression**

$$p(\mathbf{f}_{\mathcal{M}}, f_n)^{\beta_t} q_{\gamma_t}(\mathbf{f}_{\mathcal{M}}, f_n)^{1-\beta_t} := \mathcal{N} \left(\begin{bmatrix} \tilde{\mathbf{m}}_{\mathcal{M}} \\ \tilde{\mathbf{m}}_n \end{bmatrix}, \begin{bmatrix} \tilde{\mathbf{K}}_{\mathcal{M},\mathcal{M}} & \tilde{\mathbf{K}}_{\mathcal{M},n} \\ \tilde{\mathbf{K}}_{n,\mathcal{M}} & \tilde{\mathbf{K}}_{n,n} \end{bmatrix} \right) \\ \propto \mathcal{N} \left(\begin{bmatrix} 0 \\ 0 \end{bmatrix}, \begin{bmatrix} \mathbf{K}_{\mathcal{M},\mathcal{M}} & \mathbf{K}_{\mathcal{M},n} \\ \mathbf{K}_{n,\mathcal{M}} & \mathbf{K}_{n,n} \end{bmatrix} \right)^{\beta_t} \times \mathcal{N} \left(\begin{bmatrix} \mu_{\mathcal{M}} \\ \mu_n \end{bmatrix}, \begin{bmatrix} \Sigma_{\mathcal{M},\mathcal{M}} & \Sigma_{\mathcal{M},n} \\ \Sigma_{n,\mathcal{M}} & \Sigma_{n,n} \end{bmatrix} \right)^{(1-\beta_t)}$$

$$\hat{q}_{t+1}(\mathbf{f}_{\mathcal{M}}, f_n) \propto \mathcal{N}(y_n | f_n, \sigma^2 / (N\beta_t)) \times \mathcal{N} \left(\begin{bmatrix} \tilde{\mathbf{m}}_{\mathcal{M}} \\ \tilde{\mathbf{m}}_n \end{bmatrix}, \begin{bmatrix} \tilde{\mathbf{K}}_{\mathcal{M},\mathcal{M}} & \tilde{\mathbf{K}}_{\mathcal{M},n} \\ \tilde{\mathbf{K}}_{n,\mathcal{M}} & \tilde{\mathbf{K}}_{n,n} \end{bmatrix} \right)$$

- **(Nonconjugate case)** optimize an upper bound of $\text{KL}[q_{\gamma}(\mathbf{f}_{\mathcal{M}}) \|\hat{q}_{t+1}(\mathbf{f}_{\mathcal{M}})]$

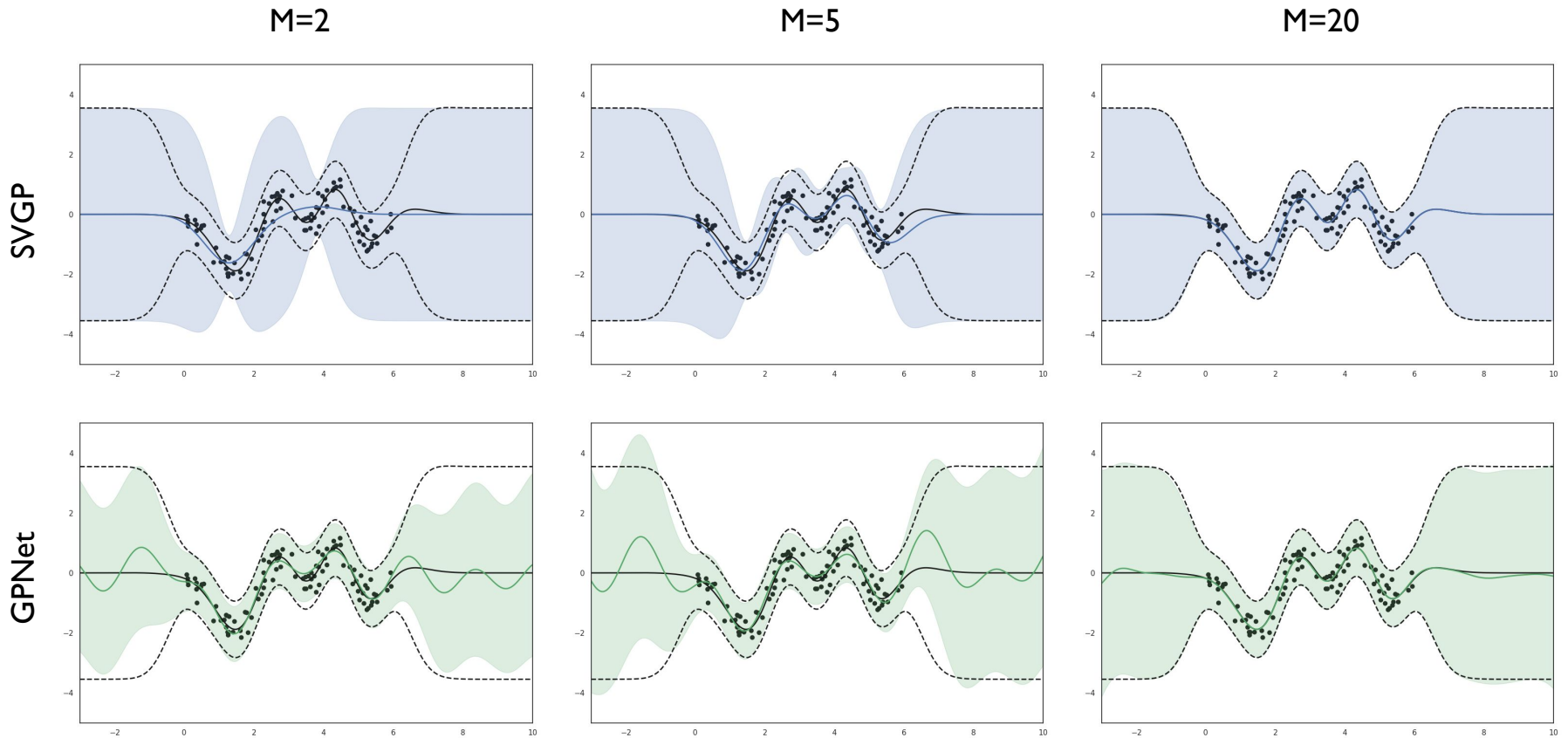
$$\min_{\gamma} \text{KL}[q_{\gamma}(\mathbf{f}_{\mathcal{M}}, f_n) \|\hat{q}_{t+1}(\mathbf{f}_{\mathcal{M}}, f_n)] \Leftrightarrow \max_{\gamma} \mathcal{L}_t(q_{\gamma}; q_{\gamma_t}, \mathbf{X}_{\mathcal{M}})$$

$$\mathcal{L}_t(q_{\gamma}; q_{\gamma_t}, \mathbf{X}_{\mathcal{M}}) = \mathbb{E}_{q_{\gamma}(\mathbf{f}_{\mathcal{M}}, f_n)} [N\beta_t \log p(y_n | f_n) + \beta_t \log p(\mathbf{f}_{\mathcal{M}}, f_n) +$$

$$(1 - \beta_t) \log q_{\gamma_t}(\mathbf{f}_{\mathcal{M}}, f_n) - \log q_{\gamma}(\mathbf{f}_{\mathcal{M}}, f_n)]$$

Scalable Training of Inference Networks for GP Models

Measurement points vs. inducing points



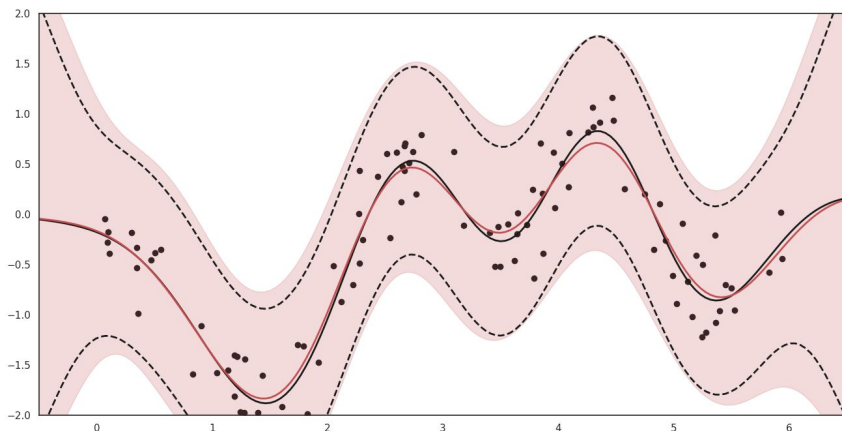
- inducing points - **expressiveness** of variational approximation
- measurement points - **variance** of training

Scalable Training of Inference Networks for GP Models

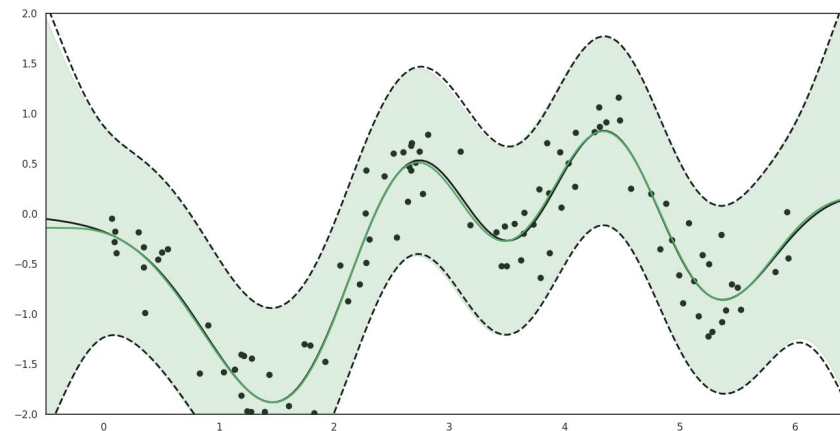
Effect of proper minibatch training

- Fix underfitting

N=100, batch size=20



FBNN, M=20



GPNet, M=20

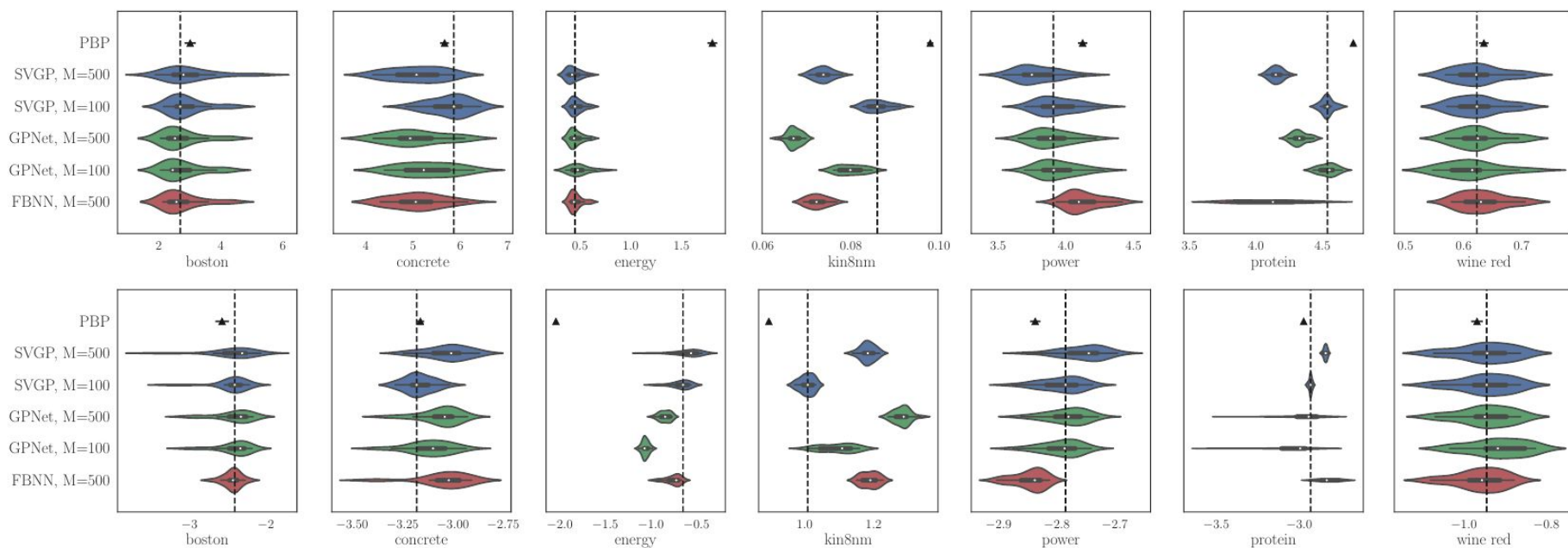
- Better performance with more measurement points

Airline Delay (700K)

METRIC	M=100			M=500		
	SVGP	GPNET	FBNN	SVGP	GPNET	FBNN
RMSE	24.261	24.055	23.801	23.698	23.675	24.114
Test LL	-4.618	-4.616	-4.586	-4.594	-4.601	-4.582

Scalable Training of Inference Networks for GP Models

Regression & Classification



Regression benchmarks

METHODS	MNIST	CIFAR10
SVGP, RBF-ARD (Krauth et al., 2016)	1.55%	-
Conv GP (van der Wilk et al., 2017)	1.22%	35.4%
SVGP, CNN-GP (Garriga-Alonso et al., 2019)	2.4%	-
GPNNet, CNN-GP	1.12%	24.63%
NN-GP (Lee et al., 2018)	1.21%	44.34%
CNN-GP (Garriga-Alonso et al., 2019)	0.96%	-
ResNet-GP (Garriga-Alonso et al., 2019)	0.84%	-
CNN-GP (Novak et al., 2019)	0.88%	32.86%

GP classification with a prior
derived from
infinite-width Bayesian ConvNets

Poster #227

Code: <https://github.com/thjashin/gp-infer-net>