INTER-DOMAIN DEEP GAUSSIAN PROCESSES ICML 2020





Tim G. J. Rudner* Dino Sejdinovic

Project website: http://bit.ly/inter-domain-dgp







University of Oxford

Yarin Gal



UNIVERSITY OF OXFORD

BAYESIAN DEEP LEARNING

Classification

Bayesian Neural Networks



Figure 1. Retina scans for diabetic retinopathy diagnosis.¹

¹ Sebastian Farquhar, Michael Osborne, Yarin Gal. Radial Bayesian Neural Networks: Beyond Discrete Support In Large-Scale Bayesian Deep Learning. In AISTATS 2020.
 ² Hugh Salimbeni, Vincent Dutordoir, James Hensman, Marc P. Deisenroth. Deep Gaussian Processes with Importance-Weighted Variational Inference. In ICML 2019.

Regression

Deep Gaussian processes (DGPs)



Figure 2. Complex, multi-modal deep GP posterior.²

BAYESIAN DEEP LEARNING

Classification

Bayesian Neural Networks



Figure 1. Retina scans for diabetic retinopathy diagnosis.¹

¹ Sebastian Farquhar, Michael Osborne, Yarin Gal. Radial Bayesian Neural Networks: Beyond Discrete Support In Large-Scale Bayesian Deep Learning. In AISTATS 2020.
 ² Hugh Salimbeni, Vincent Dutordoir, James Hensman, Marc P. Deisenroth. Deep Gaussian Processes with Importance-Weighted Variational Inference. In ICML 2019.

Regression

Deep Gaussian processes (DGPs)



Figure 2. Complex, multi-modal deep GP posterior.²

Shortcomings

- Reliant on "local" approximations based on "inducing points"
- Scales quadratically in the number of local approximations
- Does not capture any global structure in the data

² Hugh Salimbeni, Vincent Dutordoir, James Hensman, Marc P. Deisenroth. Deep Gaussian Processes with Importance-Weighted Variational Inference. In ICML 2019.

Regression

Deep Gaussian processes (DGPs)



Figure 2. Complex, multi-modal deep GP posterior.²

Shortcomings

- Reliant on "local" approximations based on "inducing points"
- Scales quadratically in the number of local approximations
- Does not capture global structure in the data



Figure 3. Multi-step function.

Shortcomings

- Reliant on "local" approximations based on "inducing points"
- Scales quadratically in the number of local approximations
- Does not capture global structure in the data



Figure 3. Multi-step function.

Shortcomings

- Reliant on "local" approximations based on "inducing points"
- Scales quadratically in the number of local approximations
- Does not capture global structure in the data



Figure 3. Multi-step function.

Shortcomings

- Reliant on "local" approximations based on "inducing points"
- Scales quadratically in the number of local approximations
- Does not capture global structure in the data

Global structure in data



Figure 5. Audio signal.

Shortcomings

- Reliant on "local" approximations based on "inducing points"
- Scales quadratically in the number of local approximations
- Does not capture global structure in the data

Global structure in data



Figure 5. Audio signal.

Shortcomings

- Reliant on "local" approximations based on "inducing points"
- Scales quadratically in the number of local approximations
- Does not capture global structure in the data

Global structure in data



Figure 5. Audio signal.

Shortcomings

- Reliant on "local" approximations based on "inducing points"
- Scales quadratically in the number of local approximations
- Does not capture global structure in the data



Shortcomings

- Reliant on "local" approximations based on "inducing points"
- Scales quadratically in the number of local approximations
- Does not capture global structure in the data



Shortcomings

- Reliant on "local" approximations based on "inducing points"
- Scales quadratically in the number of local approximations
- Does not capture global structure in the data



Shortcomings

- Reliant on "local" approximations based on "inducing points"
- Scales quadratically in the number of local approximations
- Does not capture global structure in the data



Shortcomings

- Reliant on "local" approximations based on "inducing points"
- Scales quadratically in the number of local approximations
- Does not capture global structure in the data

Solution?

Incorporate global structure into approximation



Shortcomings

- Reliant on "local" approximations based on "inducing points"
- Scales quadratically in the number of local approximations
- Does not capture global structure in the data

Solution?

Incorporate global structure into approximation



Shortcomings

- Reliant on "local" approximations based on "inducing points"
- Scales quadratically in the number of local approximations
- Does not capture global structure in the data

Solution?

Incorporate global structure into approximation



Shortcomings

- Reliant on "local" approximations based on "inducing points"
- Scales quadratically in the number of local approximations
- Does not capture global structure in the data

Solution?

Incorporate global structure into approximation



INTER-DOMAIN DEEP GAUSSIAN PROCESSES

BACKGROUND: DEEP GPS

Deep Gaussian Processes

Deep Gaussian processes with L layers

$$\mathbf{y} = \mathbf{f}^{(L)} + \boldsymbol{\epsilon} = f^{(L)}$$

- Exact inference in this model is intractable
- Requires approximate inference

 $\left(f^{(L-1)}(\dots f(\mathbf{X}))\dots\right) + \boldsymbol{\epsilon}$

BACKGROUND: INDUCING POINT APPROXIMATIONS

Variational Inference via Local Approximations

Define inducing variables

- Construct operators
- Compute posterior predictive distribution

¹ Michalis K. Titsias. Variational Learning of Inducing Variables in Sparse Gaussian Processes. In AISTATS 2009.

$u(\mathbf{z}) = f(\mathbf{z})$

$K_{fu} = K(X, Z)$ $K_{uu} \equiv K(Z, Z)$

BACKGROUND: INDUCING POINT APPROXIMATIONS

Variational Inference via Local Approximations

Define inducing variables

- Construct operators
 - $\mathbf{K}_{\mathbf{fu}} = \mathbf{K}(\mathbf{X},$
- Compute posterior predictive distribution

¹ Michalis K. Titsias. Variational Learning of Inducing Variables in Sparse Gaussian Processes. In AISTATS 2009.

$$u(\mathbf{z}) = f(\mathbf{z})$$

$$\mathbf{Z}) \quad \mathbf{K}_{\mathbf{u}\mathbf{u}} \equiv \mathbf{K}(\mathbf{Z},\mathbf{Z})$$

BACKGROUND: INDUCING POINT APPROXIMATIONS

Variational Inference via Local Approximations

Define inducing variables

Construct operators

$$\mathbf{K_{fu}} = \mathbf{K}(\mathbf{X},$$

Compute posterior predictive distribution

¹ Michalis K. Titsias. Variational Learning of Inducing Variables in Sparse Gaussian Processes. In AISTATS 2009.

$$u(\mathbf{z}) = f(\mathbf{z})$$



Local vs. Global Approximations

Local inducing variables

$$u(\mathbf{z}) = \int_{\mathbb{R}^D} f(\mathbf{x}) g(\mathbf{x},$$

$$u(\mathbf{z}) = f(\mathbf{z})$$



Local vs. Global Approximations

Local inducing variables

$$u(\mathbf{z}) = \int_{\mathbb{R}^D} f(\mathbf{x}) g(\mathbf{x},$$

$$u(\mathbf{z}) = f(\mathbf{z})$$



Local vs. Global Approximations

Local inducing variables

$$u(\mathbf{z}) = \int_{\mathbb{R}^D} f(\mathbf{x}) g(\mathbf{x},$$

$$u(\mathbf{z}) = f(\mathbf{z})$$



Local vs. Global Approximations

Local inducing variables

$$u(\mathbf{z}) = \int_{\mathbb{R}^D} f(\mathbf{x}) g(\mathbf{x},$$

$$u(\mathbf{z}) = f(\mathbf{z})$$



RKHS Fourier Features

Define truncated Fourier basis:

 $\phi(x) = [1, \cos(\omega_1(x - a)), \dots, \cos(\omega_M a)]$

Inter-domain operators

contain information about global structure4

⁴ James Hensman, Nicolas Durrande, and Arno Solin. Variational Fourier features for Gaussian processes. Journal of Machine Learning Research 2018.

 $\mathbf{K}^{\phi}_{\mathbf{fu}}$

 \mathbf{K}^{ϕ}

$$(x-a)$$
, sin $(\omega_1(x-a))$, ..., sin $(\omega_M(x-a))$] ^{\top}

$$=\phi(\mathbf{X})$$

$$= \langle \phi, \phi \rangle_{\mathcal{H}}$$

How can we use inter-domain transformations?

Damianou & Lawrence (2013)² construct posterior from:

 $\mathbf{K}_{\mathbf{f}^{\ell}\mathbf{u}^{\ell}}^{\phi^{\top}}$

² Andreas Damianou and Neil Lawrence. Deep Gaussian processes. In AISTATS 2013.

$$q\left(\mathbf{f}_{n}^{(\ell)}\right) \mathrm{d}\mathbf{f}_{n}^{(\ell)}$$

How can we use inter-domain transformations?

Damianou & Lawrence (2013)² construct posterior from:



² Andreas Damianou and Neil Lawrence. Deep Gaussian processes. In AISTATS 2013.

$$q\left(\mathbf{f}_{n}^{(\ell)}\right) \mathrm{d}\mathbf{f}_{n}^{(\ell)}$$

How can we use inter-domain transformations?

Damianou & Lawrence (2013)² construct posterior from:



² Andreas Damianou and Neil Lawrence. Deep Gaussian processes. In AISTATS 2013.

$$q\left(\mathbf{f}_{n}^{(\ell)}\right) \mathrm{d}\mathbf{f}_{n}^{(\ell)}$$

Hard to compute!

Can we do better? Yes!

- Use different factorization of variational posterior
- Doubly stochastic variational inference (DSVI)³:

$$q\left(\left\{\mathbf{F}^{(\ell)}\right\}_{\ell=1}^{L}\right) = \prod_{\ell=1}^{L} \mathcal{N}\left(\mathbf{F}^{(\ell)} | \widetilde{\mathbf{m}}^{(\ell)}, \widetilde{\mathbf{S}}^{(\ell)}\right)$$

with

$$\begin{split} \widetilde{\mathbf{m}}\left(\mathbf{F}^{(\ell)}\right) &\equiv \mathbf{m}_{\mathbf{f}^{\ell}} - \mathbf{K}_{\mathbf{f}^{\ell}\mathbf{u}^{\ell}}\mathbf{K}_{\mathbf{u}^{\ell}\mathbf{u}^{\ell}}^{-1}\left(\boldsymbol{\mu}^{(\ell)} - \mathbf{m}_{\mathbf{u}^{\ell}}\right) \\ \widetilde{\mathbf{S}}\left(\mathbf{F}^{(\ell)}, \mathbf{F}^{(\ell)}\right) &\equiv \mathbf{K}_{\mathbf{f}^{\ell}\mathbf{f}^{\ell}} - \mathbf{K}_{\mathbf{f}^{\ell}\mathbf{u}^{\ell}}\mathbf{K}_{\mathbf{u}^{\ell}\mathbf{u}^{\ell}}^{-1}\left(\mathbf{K}_{\mathbf{u}^{\ell}\mathbf{u}^{\ell}} - \boldsymbol{\Sigma}^{(\ell)}\right)\mathbf{K}_{\mathbf{u}^{\ell}\mathbf{u}^{\ell}}^{-1}\mathbf{K}_{\mathbf{u}^{\ell}\mathbf{u}^{\ell}}\right] \end{split}$$

Can we do better? Yes!

- Use different factorization of variational posterior
- Doubly stochastic variational inference (DSVI)³:

$$q\left(\left\{\mathbf{F}^{(\ell)}\right\}_{\ell=1}^{L}\right) = \prod_{\ell=1}^{L} \mathcal{N}\left(\mathbf{F}^{(\ell)} | \widetilde{\mathbf{m}}^{(\ell)}, \widetilde{\mathbf{S}}^{(\ell)}\right)$$

with

$$\begin{split} \widetilde{\mathbf{m}}\left(\mathbf{F}^{(\ell)}\right) &\equiv \mathbf{m}_{\mathbf{f}^{\ell}} - \mathbf{K}_{\mathbf{f}^{\ell}\mathbf{u}^{\ell}}\mathbf{K}_{\mathbf{u}^{\ell}\mathbf{u}^{\ell}}^{-1} \left(\boldsymbol{\mu}^{(\ell)} - \mathbf{m}_{\mathbf{u}^{\ell}}\right) \\ \widetilde{\mathbf{S}}\left(\mathbf{F}^{(\ell)}, \mathbf{F}^{(\ell)}\right) &\equiv \mathbf{K}_{\mathbf{f}^{\ell}\mathbf{f}^{\ell}} - \mathbf{K}_{\mathbf{f}^{\ell}\mathbf{u}^{\ell}}\mathbf{K}_{\mathbf{u}^{\ell}\mathbf{u}^{\ell}}^{-1} \left(\mathbf{K}_{\mathbf{u}^{\ell}\mathbf{u}^{\ell}} - \boldsymbol{\Sigma}^{(\ell)}\right) \mathbf{K}_{\mathbf{u}^{\ell}\mathbf{u}^{\ell}}^{-1} \mathbf{K}_{\mathbf{u}^{\ell}\mathbf{f}^{\ell}} \end{split}$$

Inter-domain Deep GP approximate posterior

- Use inter-domain operators $\mathbf{K}^{\phi}_{\mathbf{f}_{\mathbf{H}}}$ and $\mathbf{K}^{\phi}_{\mathbf{H}_{\mathbf{H}}}$ as drop-in replacements
- Posterior predictive distribution under DSVI³:

$$q\left(\left\{\mathbf{F}^{(\ell)}\right\}_{\ell=1}^{L}\right) = \prod_{\ell=1}^{L} \mathcal{N}\left(\mathbf{F}^{(\ell)} | \widetilde{\mathbf{m}}^{(\ell)}, \widetilde{\mathbf{S}}^{(\ell)}\right)$$

with

$$\begin{split} \widetilde{\mathbf{m}}\left(\mathbf{F}^{(\ell)}\right) &\equiv \mathbf{m}_{\mathbf{f}^{\ell}} - \mathbf{K}_{\mathbf{f}'\mathbf{u}^{\ell}}^{\phi} \mathbf{K}_{\mathbf{u}^{\ell}\mathbf{u}^{\ell}}^{\phi^{-1}} \left(\boldsymbol{\mu}^{(\ell)} - \mathbf{m}_{\mathbf{u}^{\ell}}\right) \\ \widetilde{\mathbf{S}}\left(\mathbf{F}^{(\ell)}, \mathbf{F}^{(\ell)}\right) &\equiv \mathbf{K}_{\mathbf{f}^{\ell}\mathbf{f}^{\ell}} - \mathbf{K}_{\mathbf{f}^{\ell}\mathbf{u}^{\ell}}^{\phi} \mathbf{K}_{\mathbf{u}^{\ell}\mathbf{u}^{\ell}}^{\phi^{-1}} \left(\mathbf{K}_{\mathbf{u}^{\ell}\mathbf{u}^{\ell}}^{\phi} - \boldsymbol{\Sigma}^{(\ell)}\right) \mathbf{K}_{\mathbf{u}^{\ell}\mathbf{u}^{\ell}}^{\phi^{-1}} \mathbf{K}_{\mathbf{u}^{\ell}\mathbf{f}^{\ell}}^{\phi} \end{split}$$

Inter-domain Deep GP approximate posterior

- Use inter-domain operators $\mathbf{K}^{\phi}_{\mathbf{fu}}$ and $\mathbf{K}^{\phi}_{\mathbf{uu}}$ as drop-in replacements
- Posterior predictive distribution under DSVI³:

$$q\left(\left\{\mathbf{F}^{(\ell)}\right\}_{\ell=1}^{L}\right) = \prod_{\ell=1}^{L} \mathcal{N}\left(\mathbf{F}^{(\ell)} | \widetilde{\mathbf{m}}^{(\ell)}, \widetilde{\mathbf{S}}^{(\ell)}\right)$$

with

$$\widetilde{\mathbf{m}}\left(\mathbf{F}^{(\ell)}\right) \equiv \mathbf{m}_{\mathbf{f}^{\ell}} - \mathbf{K}_{\mathbf{f}'\mathbf{u}^{\ell}}^{\phi}\mathbf{K}_{\mathbf{u}^{\ell}\mathbf{u}^{\ell}}^{\phi^{-1}}\left(\boldsymbol{\mu}^{(\ell)} - \mathbf{m}_{\mathbf{u}^{\ell}}\right)$$
$$\widetilde{\mathbf{S}}\left(\mathbf{F}^{(\ell)}, \mathbf{F}^{(\ell)}\right) \equiv \mathbf{K}_{\mathbf{f}^{\ell}\mathbf{f}^{\ell}} - \mathbf{K}_{\mathbf{f}^{\ell}\mathbf{u}^{\ell}}^{\phi}\mathbf{K}_{\mathbf{u}^{\ell}\mathbf{u}^{\ell}}^{\phi^{-1}}\left(\mathbf{K}_{\mathbf{u}^{\ell}\mathbf{u}^{\ell}}^{\phi} - \boldsymbol{\Sigma}^{(\ell)}\right)\mathbf{K}_{\mathbf{u}^{\ell}\mathbf{u}^{\ell}}^{\phi^{-1}}\mathbf{K}_{\mathbf{u}^{\ell}\mathbf{f}^{\ell}}^{\phi}$$

Inter-domain Deep GP approximate posterior

- Use inter-domain operators $\mathbf{K}^{\phi}_{\mathbf{fu}}$ and $\mathbf{K}^{\phi}_{\mathbf{uu}}$ as drop-in replacements
- Posterior predictive distribution under DSVI³:

$$q\left(\left\{\mathbf{F}^{(\ell)}\right\}_{\ell=1}^{L}\right) =$$

with

$$\widetilde{\mathbf{m}}\left(\mathbf{F}^{(\ell)}
ight) \equiv \mathbf{m}_{\mathbf{f}^{\ell}} \widetilde{\mathbf{S}}\left(\mathbf{F}^{(\ell)}, \mathbf{F}^{(\ell)}
ight) \equiv \mathbf{K}_{\mathbf{f}^{\ell}\mathbf{f}^{\ell}} - \mathbf{K}_{\mathbf{f}^{\ell}}^{\phi}$$





EMPIRICAL EVALUATION





















whereas the local approximation (left) is not.





whereas the local approximation (left) is not.





whereas the local approximation (left) is not.







			1			
200	400	600	800	1000	1200	1400





































































SUMMARY

- Exploit global structure in the data
- Better predictive performance
- Higher computational efficiency
- Simple drop-in replacement



SUMMARY

- Exploit global structure in the data
- Better predictive performance
- Higher computational efficiency
- Simple drop-in replacement



- Exploit global structure in the data
- Better predictive performance
- Higher computational efficiency
- Simple drop-in replacement



- Exploit global structure in the data
- Better predictive performance
- Higher computational efficiency
- Simple drop-in replacement

- Exploit global structure in the data
- Better predictive performance
- Higher computational efficiency
- Simple drop-in replacement

THANK YOU! Project Website:

http://bit.ly/inter-domain-dgp

EMAIL:	tim.r
WEBSITE:	http:
TWITTER:	@TIMF

udner@cs.ox.ac.uk

//timrudner.com

RUDNER