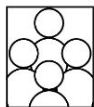


Projected Bayesian Neural Networks:

Avoiding weight-space pathologies by
learning latent representations of neural network weights

ACML, Nov 17th, 2019

Melanie F. Pradier



CRCS Center for Research on
Computation and Society

at Harvard John A. Paulson School of Engineering and Applied Sciences



HDSI | Harvard Data
Science Initiative

Projected Bayesian Neural Networks:

Avoiding weight-space pathologies by
learning latent representations of neural network weights

Joint work with my collaborators...



Weiwei Pan



Jiayu Yao



Soumya Ghosh



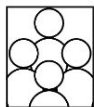
Finale Doshi-Velez

Projected Bayesian Neural Networks:

Avoiding weight-space pathologies by
learning latent representations of neural network weights

Nov 17th, 2019

Melanie F. Pradier



CRCS Center for Research on
Computation and Society

at Harvard John A. Paulson School of Engineering and Applied Sciences



HDSI | Harvard Data
Science Initiative

Success of Deep Learning



Success of Deep Learning



Success of Deep Learning



Success of Deep Learning



Black Jeans

Blue Dress

Blue Jeans



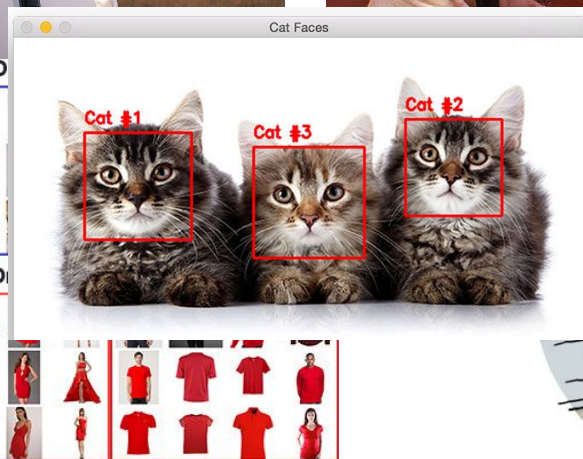
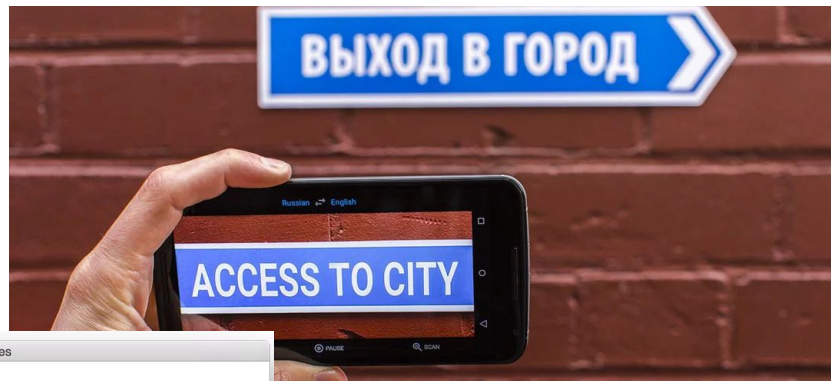
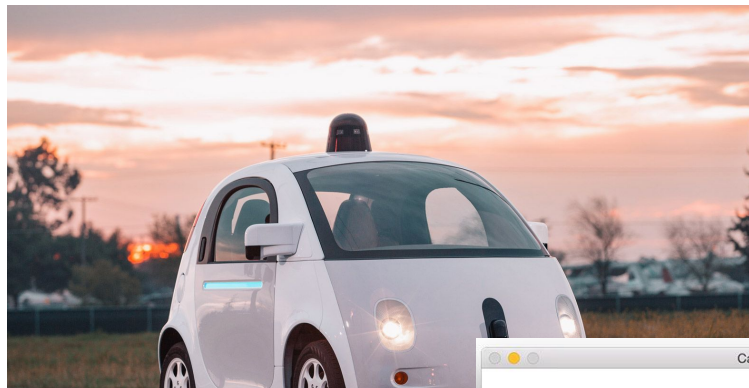
Blue Shirt

Red Dress

Red Shirt



Huge amount of opportunities...

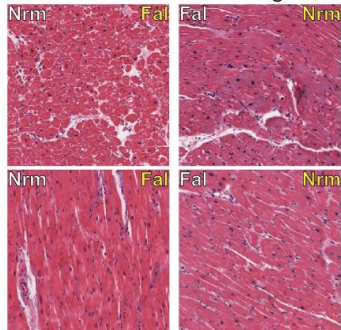


...but careful in high-stake decisions!

Deep Learning errors

False Positives False Negatives

[Nirschi et.al, 2018]



[Eykholt et.al, 2018]

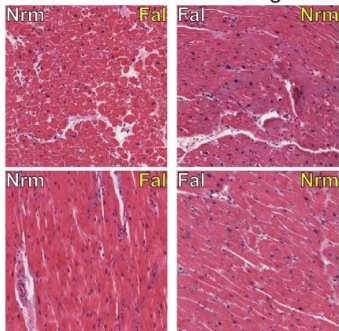


...but careful in high-stake decisions!

Deep Learning errors

False Positives False Negatives

[Nirschi et.al, 2018]



[Eykholt et.al, 2018]

Our Goal: Quantify Uncertainty

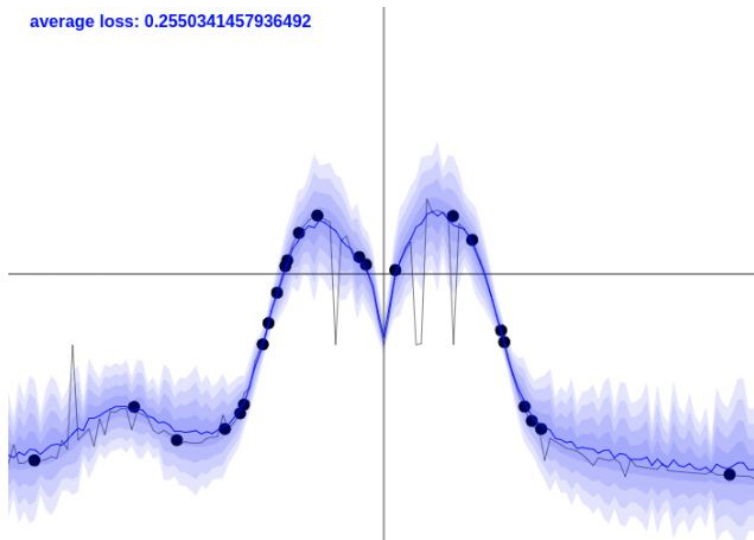
With such uncertainty, we can:

- Alert humans in unclear situations
- Diagnose ML systems (when and how does it fail)
- Get better predictive accuracy

Focus: Bayesian Neural Networks



Bayesian Neural Network (BNN)



[What my deep model does not know, post of Yarin Gal, 2015]

$$\mathcal{D} = \{\mathbf{x}_i, \mathbf{y}_i\}_{i=1}^N$$

$$\mathbf{y} = f_{\mathbf{w}}(\mathbf{x}) + \epsilon$$

$$\mathbf{w} \sim \mathcal{N}(0, \sigma_w^2 \mathbf{I}), \quad \epsilon \sim \mathcal{N}(0, \sigma_\epsilon^2 \mathbf{I})$$

Quantities of interest:

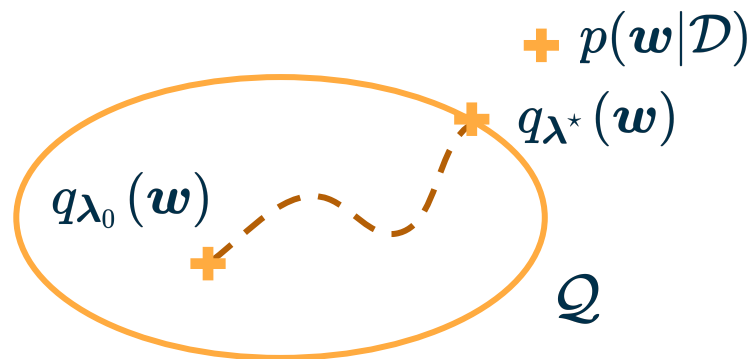
- Posterior of the weights $p(\mathbf{w}|\mathcal{D})$
- Predictive distribution

$$p(\mathbf{y}^*|\mathbf{x}^*, \mathcal{D}) = \int p(\mathbf{y}^*|\mathbf{x}^*, \mathbf{w})p(\mathbf{w}|\mathcal{D})d\mathbf{w}$$

Key challenge: $p(\boldsymbol{w}|\mathcal{D})$ intractable

Variational Inference: appealing for scalability

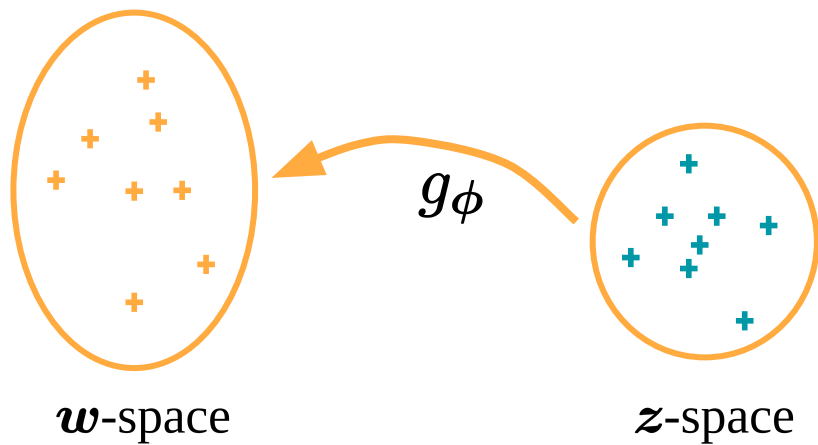
- Mean-field Approx. [Graves, 1993] [Blundell et.al, 2015]
- Structured Variational Approximations
 - Multivariate Gaussians [Louizos et.al, 2016; Sun et.al, 2017]
 - Hierarchical Variational Models [Ranganath et.al, 2016]
- Normalizing Flows and Transformations
 - Multiplicative Normalizing Flow [Louizos et. al, 2017]
 - Hypernetworks [Krueger et.al, 2017; Pawłowski et.al, 2017]



$$q_{\lambda}(\boldsymbol{w}) \in \mathcal{Q}$$

Our Approach

Projected Bayesian NN (proj-BNN)



$$D_w \gg D_z$$

$$\mathcal{D} = \{\mathbf{x}_i, \mathbf{y}_i\}_{i=1}^N$$

$$\mathbf{y} = f_w(\mathbf{x}) + \epsilon \quad \epsilon \sim \mathcal{N}(0, \sigma_\epsilon^2 \mathbf{I})$$

$$\mathbf{w} = g_\phi(\mathbf{z}), \quad \mathbf{z} \sim p(\mathbf{z}), \quad \phi \sim p(\phi),$$

Quantities of interest:

- Posterior of the weights $p(\mathbf{z}, \phi | \mathcal{D})$
- Predictive distribution

$$p(\mathbf{y}^* | \mathbf{x}^*, \mathcal{D}) = \int p(\mathbf{y}^* | \mathbf{x}^*, \mathbf{z}, \phi) p(\mathbf{z}, \phi | \mathcal{D}) d\mathbf{w}$$

How about inference?

Objective: approximate $p(\mathbf{w}|\mathcal{D})$

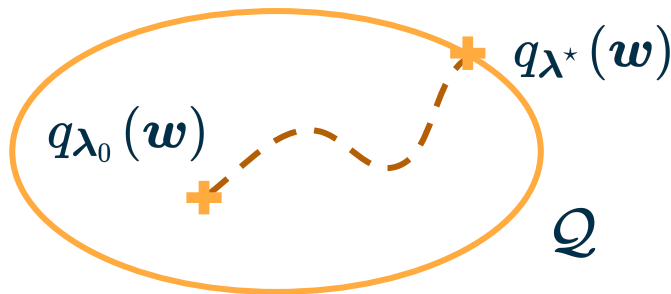
+ $p(\mathbf{w}|\mathcal{D})$

$$q_{\lambda}(\mathbf{w}) \in \mathcal{Q}$$

$$\operatorname{argmin}_{\lambda^*} D_{\text{KL}}(q_{\lambda}(\mathbf{w})||p(\mathbf{w}|\mathcal{D}))$$



$$\operatorname{argmax}_{\lambda^*} \mathcal{L}(\lambda) = \mathbb{E}_q \left[\log p(\mathbf{y}|\mathbf{x}, \mathbf{w}) \right] - D_{\text{KL}}(q_{\lambda}(\mathbf{w})||p(\mathbf{w}))$$



How about inference?

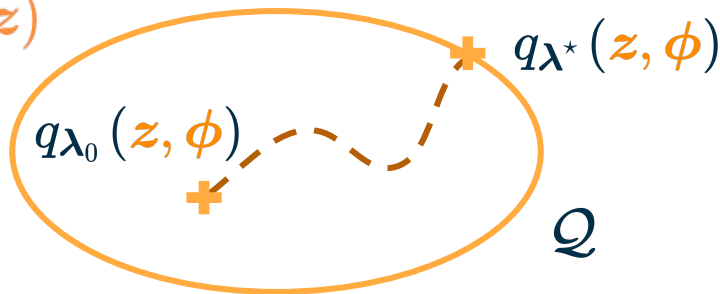
Objective: approximate $p(\mathbf{z}, \phi | \mathcal{D})$ + $p(\mathbf{z}, \phi | \mathcal{D})$

$$\mathbf{z} \sim q_{\lambda_z}(\mathbf{z}), \quad \phi \sim q_{\lambda_\phi}(\phi), \quad \mathbf{w} = g_\phi(\mathbf{z})$$

$$\operatorname{argmin}_{\lambda^*} D_{\text{KL}}(q_{\lambda}(\mathbf{z}, \phi) \| p(\mathbf{z}, \phi | \mathcal{D}))$$

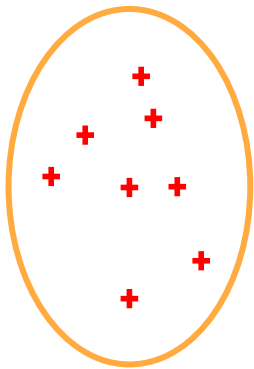


$$\operatorname{argmax}_{\lambda^*} \mathcal{L}(\boldsymbol{\lambda}) = \mathbb{E}_q \left[\log p(\mathbf{y} | \mathbf{x}, g_\phi(\mathbf{z})) \right] - D_{\text{KL}}(q_{\lambda_z}(\mathbf{z}) \| p(\mathbf{z})) - D_{\text{KL}}(q_{\lambda_\phi}(\phi) \| p(\phi))$$



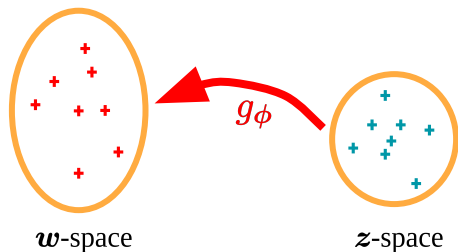
Smart-initialization Procedure

1. Characterize weight space



Sample multiple weight sets [Izmailov et.al, 2018]

2. Find point estimate g_ϕ



Train an autoencoder

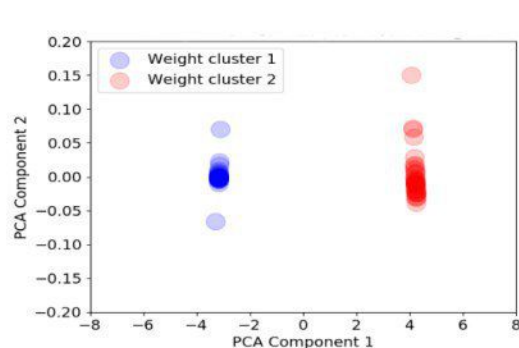
3. Black-box VI (BBVI)

$$D_{\text{KL}} \left(q_\lambda(z, \phi) \parallel p(z, \phi | \mathcal{D}) \right)$$

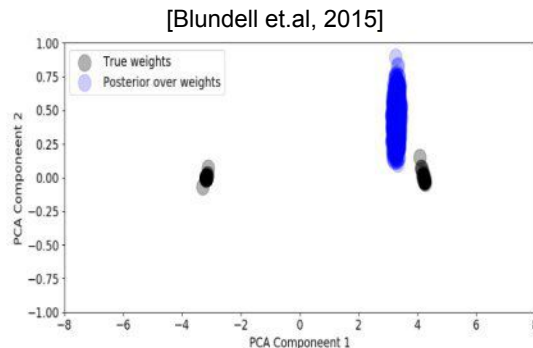
BBVI with smart initialization g_ϕ

Results

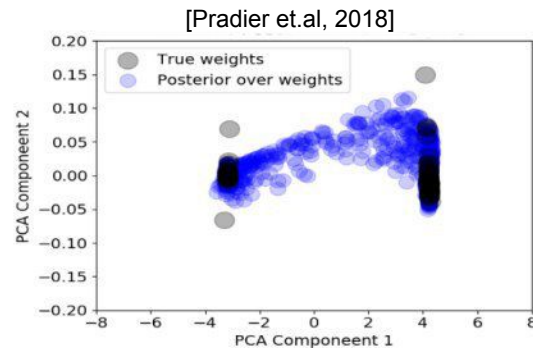
Results: Illustrative Toy Example



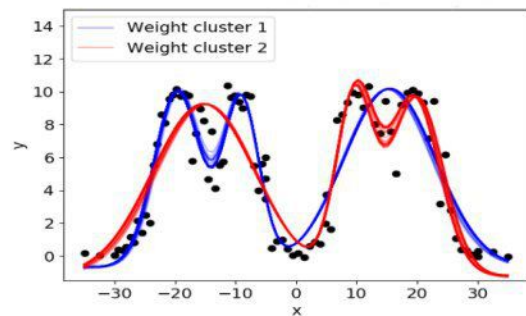
(a) Projection of true weights



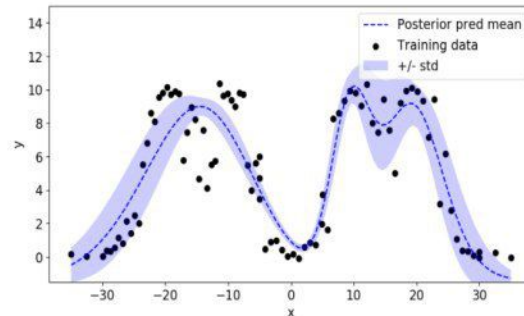
(b) BbB posterior over weights



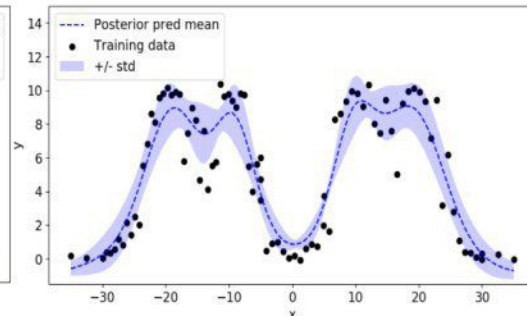
(c) Proj-BNN posterior over weights



(d) Functions from true weights

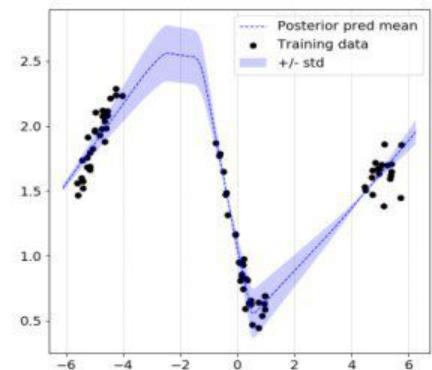


(e) BbB posterior predictive

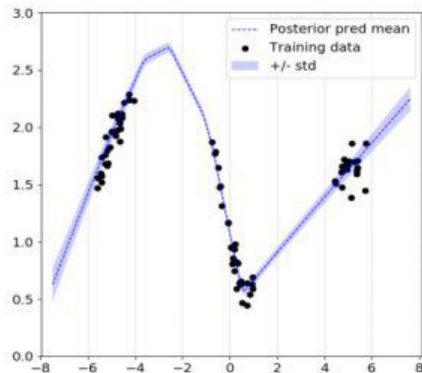


(f) Proj-BNN posterior predictive

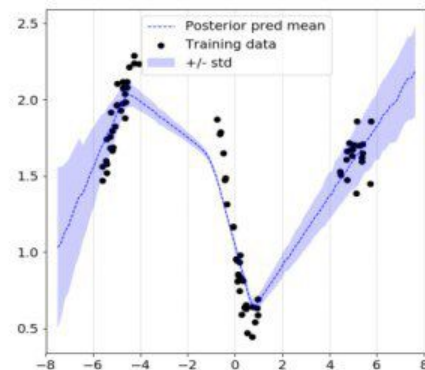
Results: Uncertainty estimation



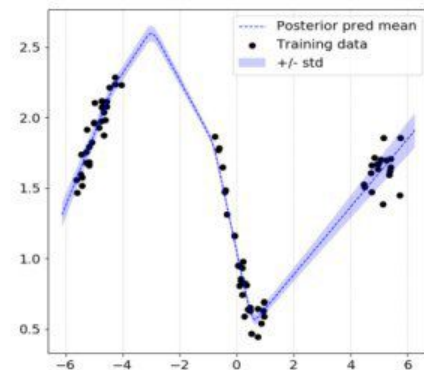
(a) Proj-BNN ($D_z = 2$)



(b) BbB



(c) MNF

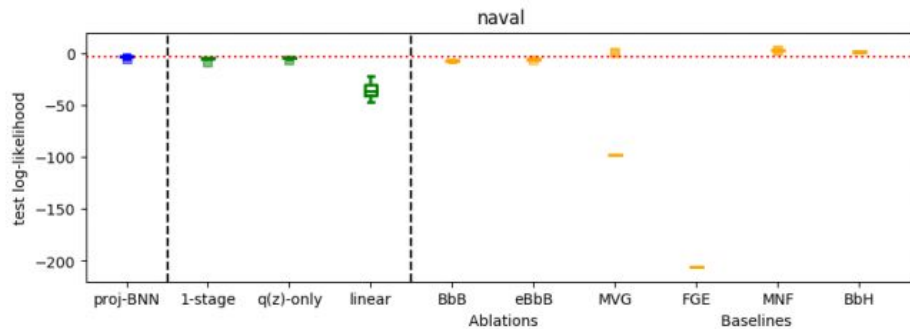
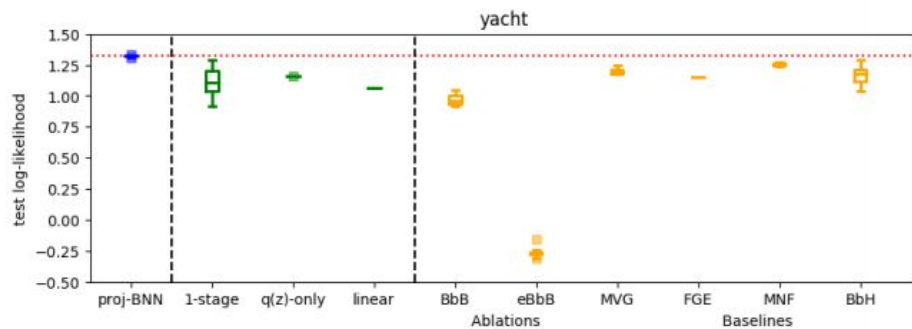
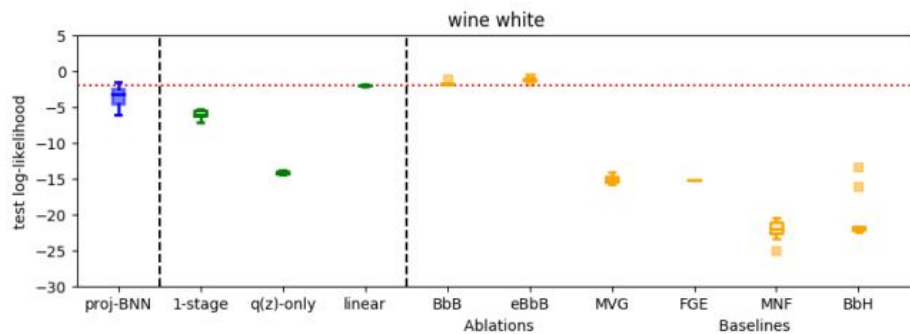
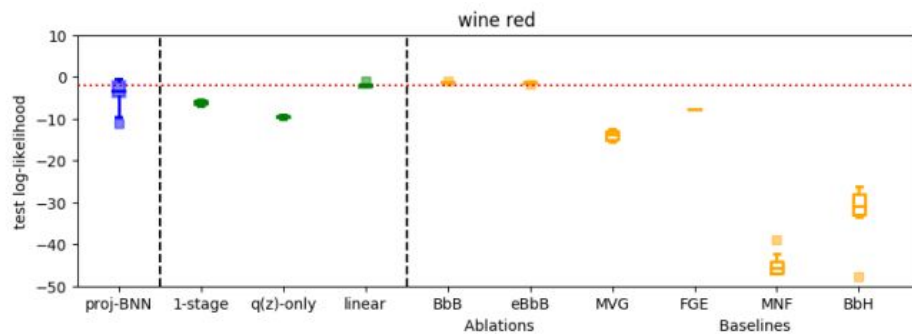


(d) MVG

- BbB: Bayes by Back Prop [Blundell et.al, 2015]
- MVG: Multivariate Gaussians [Louizos et.al, 2016]
- MNF: Multiplicative Normalizing Flow [Louizos et. al, 2017]

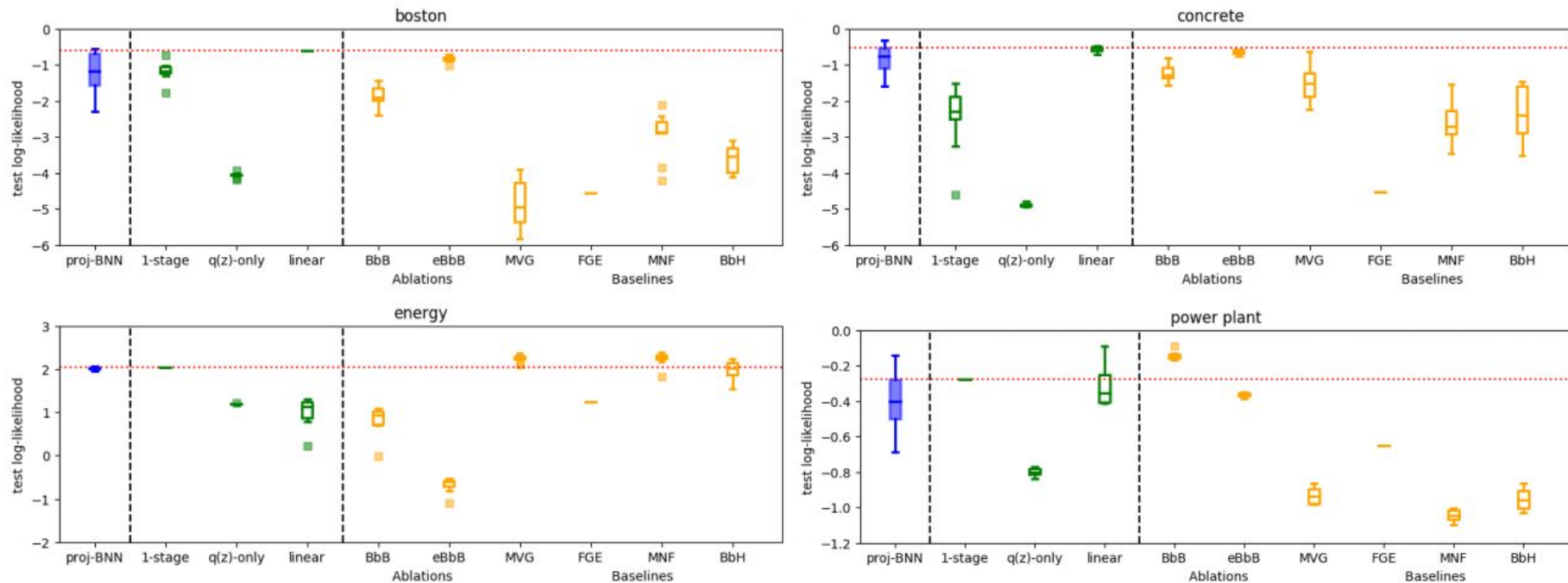
Results: Generalization

<https://arxiv.org/abs/1811.07006>



Results: Generalization

<https://arxiv.org/abs/1811.07006>



Conclusions

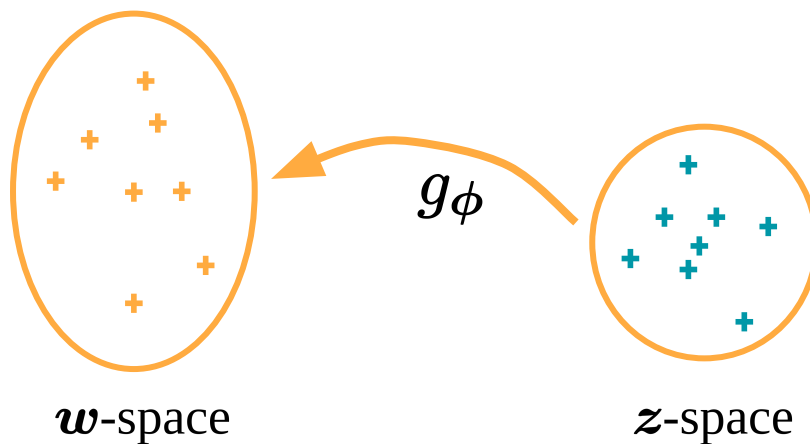


Contact: melanie@seas.harvard.edu

<https://melaniefp.github.io/>

In this talk...

- Alternative modeling for BNNs
- Better approximate inference



<https://arxiv.org/abs/1811.07006>

Thank you for listening!

Open questions

- Better evaluation of uncertainty?

“Test log likelihood can be misleading”

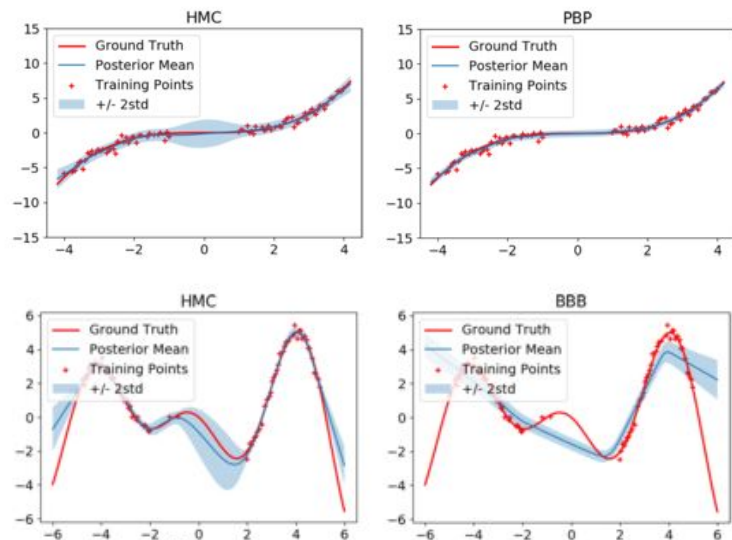
> Entangled sources of error: model, variational approx, optimization

- How does the topology in weight space looks like?

> Intuition misleading in high dimensions!

- How to exploit latent structure for interpretability?

[Yao et. al, ICML Workshop, 2019]

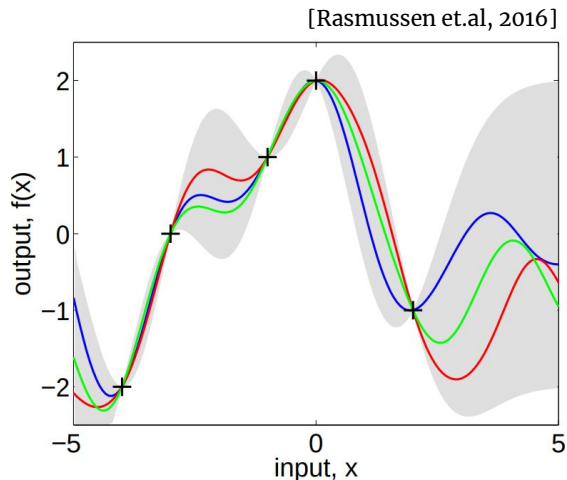


Related works on weight embeddings [Karaletsos et.al, 2018; Izmailov et.al, 2019]

Appendix

Uncertainty Estimation via GPs?

Gaussian Process (GP)



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

Drawbacks of GPs

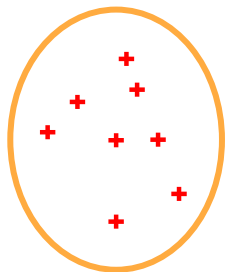
- Scalability
- Kernel learning is not trivial

Alternative: Neural Networks with uncertainty

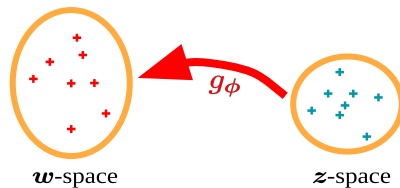
- Ensemble of Neural Networks
[Lakshminarayanan et al., 2017; Pearce et.al, 2018]
- Bayesian Neural Networks
[Buntine et al., 1991; MacKay, 1992; Neal, 1993]

Results: Generalization (Ablations)

1. Characterize w-space






2. Find point estimate g_ϕ



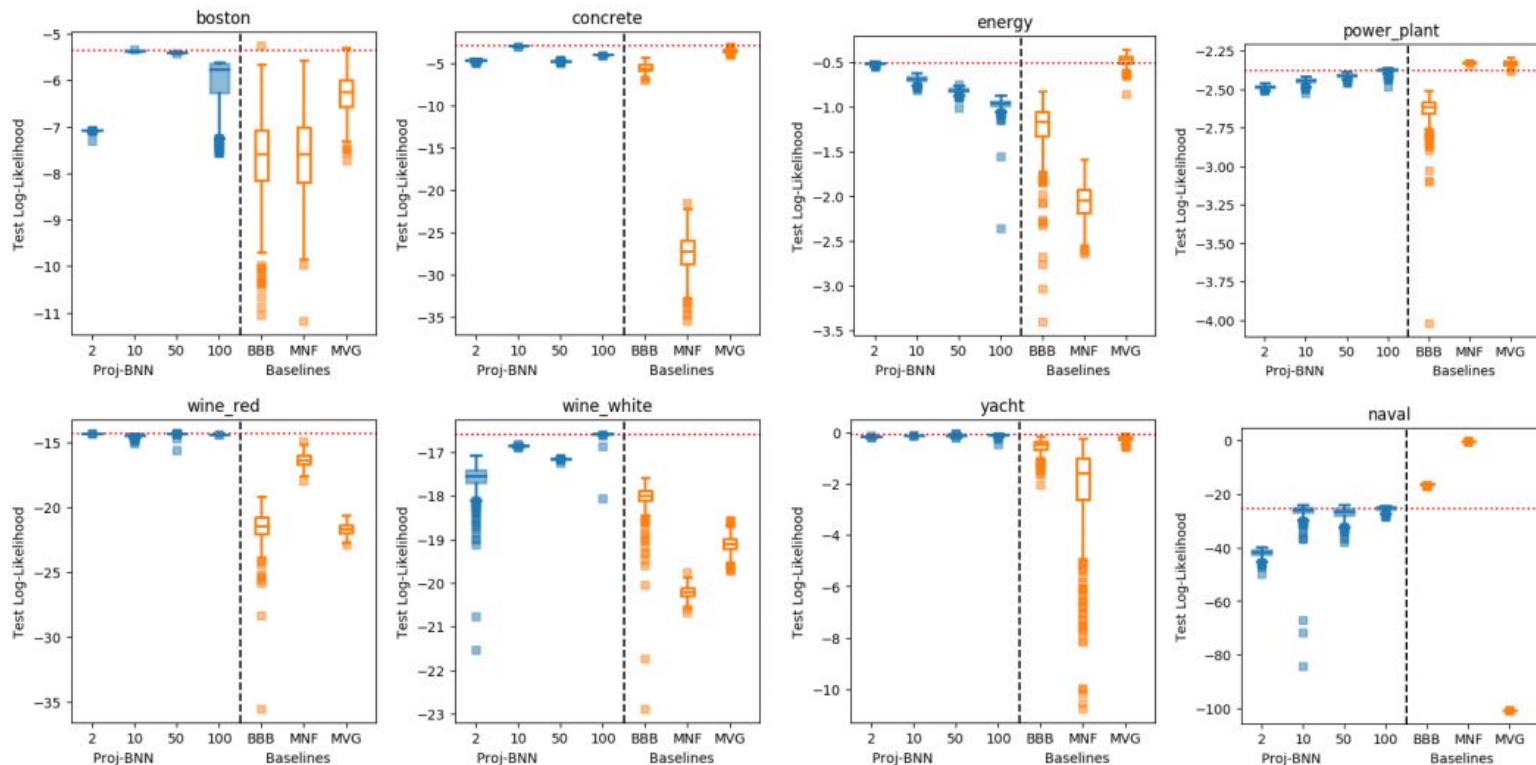
3. Black-box VI (BBVI)

$$D_{\text{KL}} \left(q_\lambda(z, \phi) || p(z, \phi | \mathcal{D}) \right)$$

1-stage			
linear		linear	
q(z) only			$q_{\lambda_z}(z)$

Cross-validation of latent dimension

<https://arxiv.org/abs/1811.07006>



Prediction-constrained Autoencoder

$$\{\boldsymbol{\theta}^*, \phi^*\} = \operatorname{argmin}_{\boldsymbol{\theta}, \phi} \mathcal{L}(\boldsymbol{\theta}, \phi) = \min_{\boldsymbol{\theta}, \phi} \left\{ \frac{1}{R} \sum_{r=1}^R \left(\mathbf{w}_{\mathbf{c}}^{(r)} - g_{\phi} \left(f_{\boldsymbol{\theta}} \left(\mathbf{w}_{\mathbf{c}}^{(r)} \right) \right) + \gamma^{(r)} \right)^2 \right. \\ \left. + \beta \mathbb{E}_{(x,y) \sim \mathcal{D}} \left[\frac{1}{R} \sum_{r=1}^R \log p(y|x, g_{\phi} \left(f_{\boldsymbol{\theta}} \left(\mathbf{w}_{\mathbf{c}}^{(r)} \right) \right)) \right] \right\},$$

My research: probabilistic models for societal needs

Highly driven by real-world application, with special emphasis on...

A) Latent Representation Learning

M. F. Pradier, B. Reis, L. Jukofsky, F. Milletti, T. Ohtomo, F. Perez-Cruz, and O. Puig. **Case-control Indian Buffet Process identifies biomarkers of response to Codrituzumab**. *BMC Cancer*. 2019.

I. Valera, M. F. Pradier, M. Lomeli, and Z. Ghahramani. **General Latent Feature Models for Heterogeneous Datasets**. *In submission to Journal of Machine Learning Research*. 2018.

M. F. Pradier, W. Pan, M. Yau, R. Singh, and F. Doshi-Velez. **Hierarchical Stick-breaking Paintbox**. *BNP@NeurIPS Workshop*. Montreal (Canada), December 2018.

B) Uncertainty Quantification

M. F. Pradier, W. Pan, J. Yao, S. Ghosh, and F. Doshi-Velez. **Projected BNNs: Avoiding Pathologies in Weight Space by projecting Neural Network Weights**. Arxiv. 2019.

B. Coker, M. F. Pradier, and F. Doshi-Velez. **Poisson Process Radial Basis Function Networks**. (Arxiv coming soon)

W. Yang, L. Lorch, M. A. Graule, S. Srinivasan, A. Suresh, J. Yao, M. F. Pradier, and F. Doshi-Velez. **Output-Constrained Bayesian Neural Networks**. *ICML Workshop on Generalization*. 2019.