



LATENT PROJECTION BNNS: AVOIDING WEIGHT-SPACE PATHOLOGIES BY LEARNING LATENT REPRESENTATIONS OF NEURAL NETWORK WEIGHTS

Melanie F. Pradier¹

Weiwei Pan¹

Jiayu Yao¹

melanie@seas.harvard.edu

Soumya Ghosh²

Finale Doshi-Velez¹

¹Harvard University

²IBM Research

INTRODUCTION

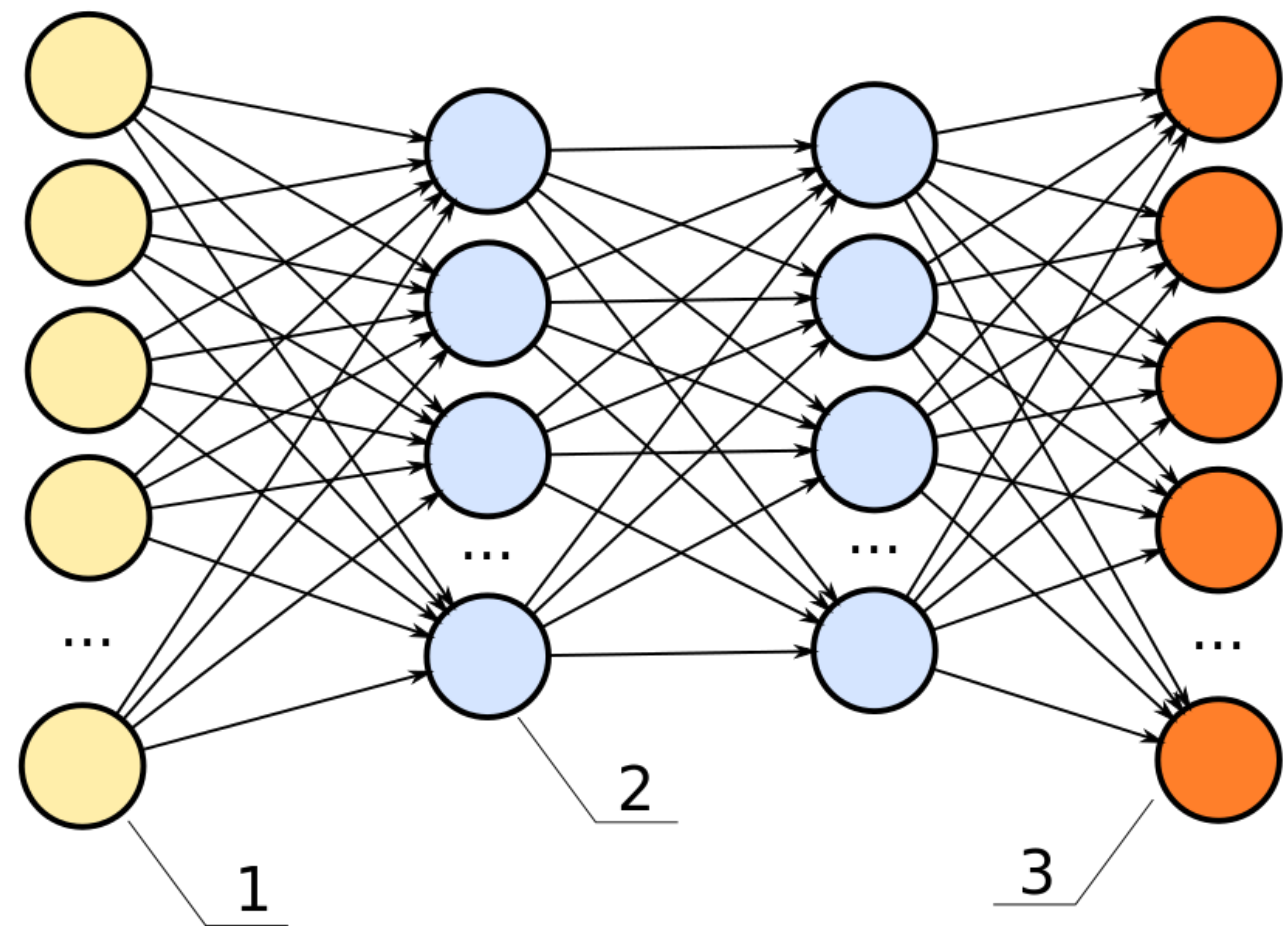
- ▶ Deep neural networks provide high-predictive accuracy BUT:
 - ▶ tend to overfit with small number of samples
 - ▶ do not provide uncertainties
- ▶ Bayesian neural networks (BNN) provide uncertainty in network weights, but inference is hard:
 - ▶ high-dimensionality of network parameter space
 - ▶ correlation between these parameters
- ▶ **Our contribution:** novel inference framework for BNNs
 - a) we capture complex distributions in high-dim weight space via low-dimensional latent space
 - b) we do inference in low-dimensional representation
 - c) better estimation of uncertainty

BAYESIAN NEURAL NETWORK (BNN)

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

$$\mathbf{w} \sim p(\mathbf{w}),$$

where $\epsilon \sim \mathcal{N}(0, \sigma_y^2 \mathbf{I})$ and typically, $w_{ij} \sim \mathcal{N}(0, \sigma_w^2)$.



1. input layer
2. hidden nodes
3. output layer

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

How can we approximate posterior distribution $p(\mathbf{w}|\mathcal{D})$ of weights \mathbf{w} ?

VARIATIONAL INFERENCE

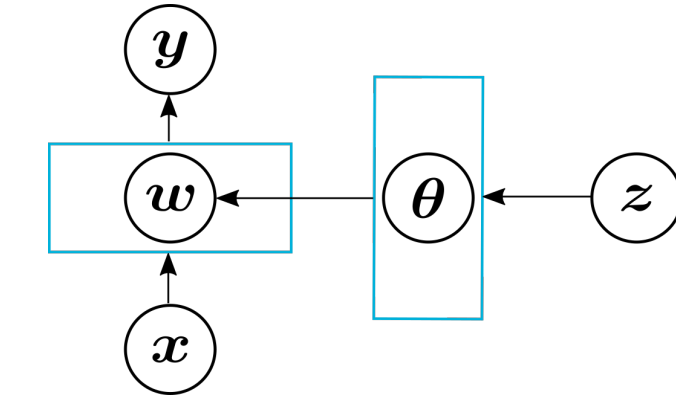
- ▶ Bayesian inference becomes an optimization problem
- ▶ Idea: choose variational distribution $q_{\lambda}(\mathbf{w})$ and minimize $D_{\text{KL}}(q_{\lambda}(\mathbf{w})||p(\mathbf{w}|\mathcal{D}))$.
- ▶ Minimizing KL divergence is equivalent to maximizing an upper bound $\mathcal{L}(\lambda)$ on the marginal likelihood of the data.

$$D_{\text{KL}}(q_{\lambda}(\mathbf{w})||p(\mathbf{w}|\mathcal{D})) = q_{\lambda}(\mathbf{w})[\log q_{\lambda}(\mathbf{w}) - \log p(\mathbf{w}|\mathcal{D})]$$

$$= -\mathcal{H}(q) - \mathbb{E}_q[\log p(\mathcal{D}, \mathbf{w})] + \log p(\mathcal{D})$$

$$= -\mathcal{L}(\lambda) + \log p(\mathcal{D})$$

LATENT PROJECTION BNN



- ▶ \mathbf{w} : network weights
- ▶ \mathbf{z} : low-dimensional latent representation
- ▶ ϕ : params of non-linear projection $g_{\phi}(\cdot)$

- ▶ We choose the following variational distribution $q_{\lambda}(\mathbf{w})$:

$$\mathbf{z} \sim q_{\lambda_z}(\mathbf{z})$$

$$\phi \sim q_{\lambda_{\phi}}(\phi)$$

$$\mathbf{w} = g_{\phi}(\mathbf{z})$$

- ▶ Novel inference framework:

1. **Characterize the space of plausible weights:** we train an ensemble of (non-Bayesian) neural networks from R multiple restarts $\rightarrow \mathbf{w}_c$

2. **Learn projection parameters ϕ :**

we train an autoencoder with a prediction-constrained loss function:

$$\min_{\theta, \phi} \left\{ \frac{1}{R} (\mathbf{w}_c - \hat{\mathbf{w}}_c + \gamma)^2 + \beta \mathbb{E}_{(x,y) \sim \mathcal{D}} [\log p(y|x, \hat{\mathbf{w}}_c)] \right\},$$

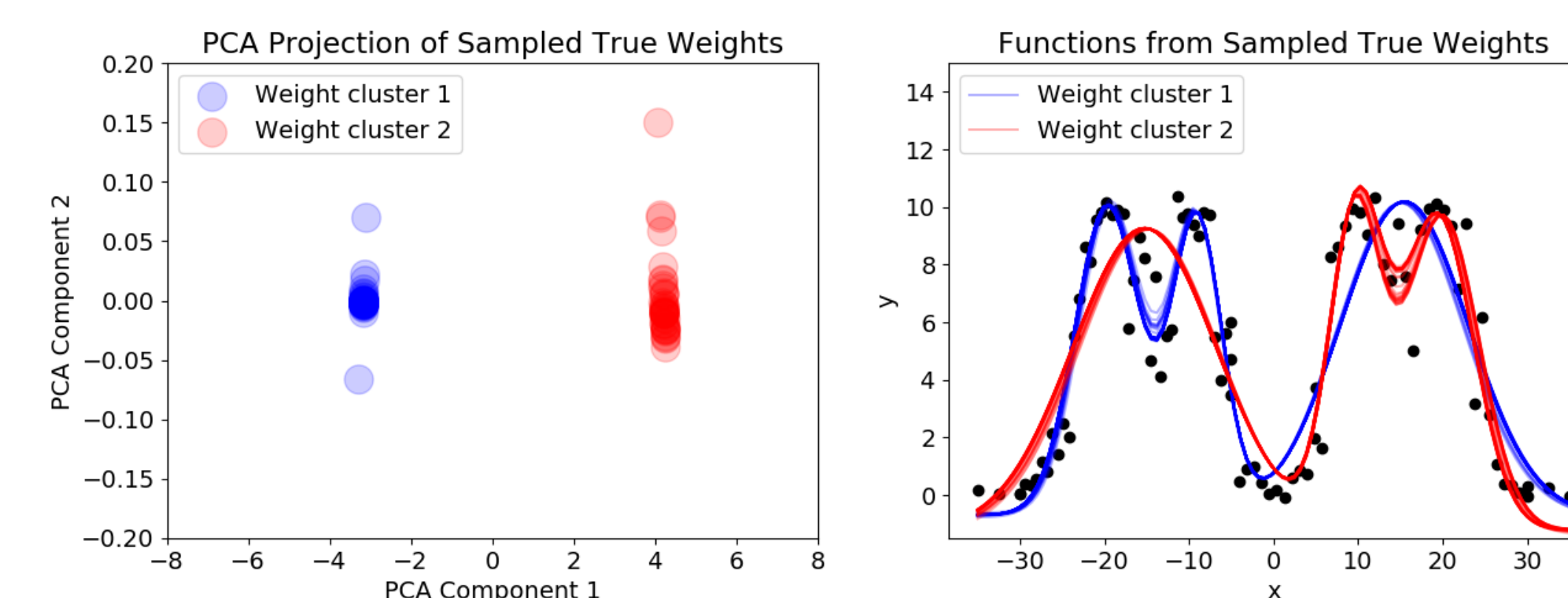
where $f_{\theta}(\cdot)$ is the encoder, $g_{\phi}(\cdot)$ is the decoder, and $\hat{\mathbf{w}}_c = g_{\phi}(f_{\theta}(\mathbf{w}_c))$.

3. **Learn variational distribution $q(\mathbf{z}, \phi)$:**

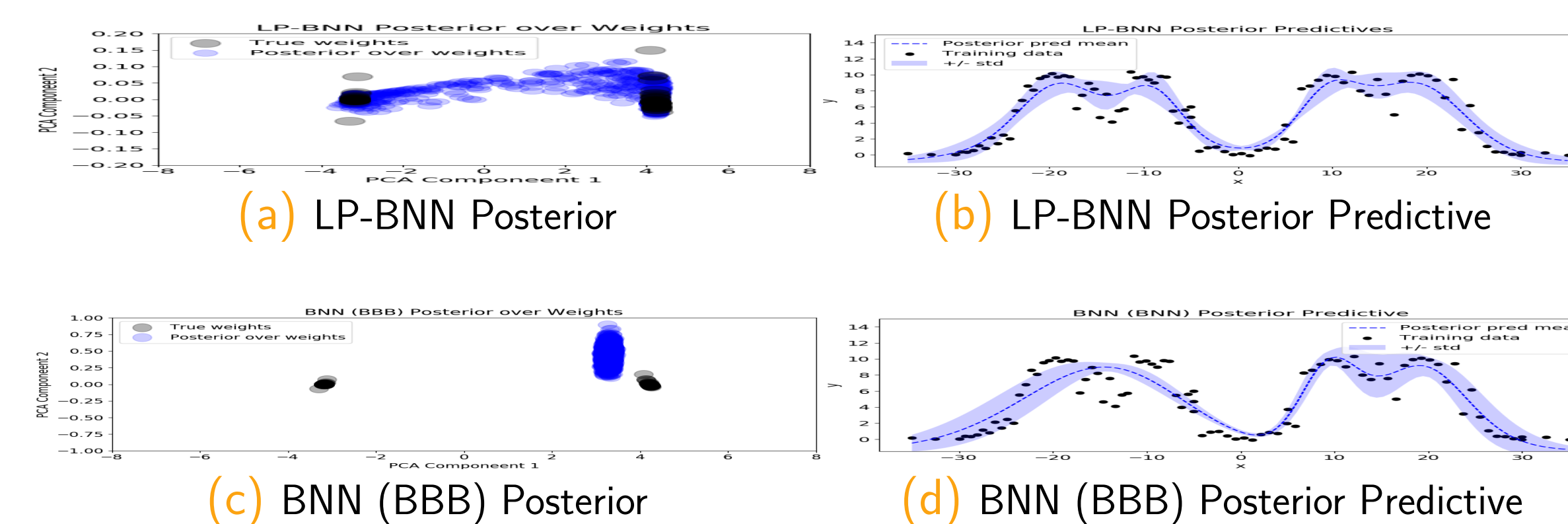
we perform BBVI with local reparametrization trick to maximize $\mathcal{L}(\lambda)$:

$$\mathcal{L}(\lambda) = \mathbb{E}_q [\log p(\mathbf{y}|\mathbf{x}, g_{\phi}(\mathbf{z}))] - D_{\text{KL}}(q_{\lambda_z}(\mathbf{z})||p(\mathbf{z})) - D_{\text{KL}}(q_{\lambda_{\phi}}(\phi)||p(\phi)).$$

PROOF-OF-CONCEPT

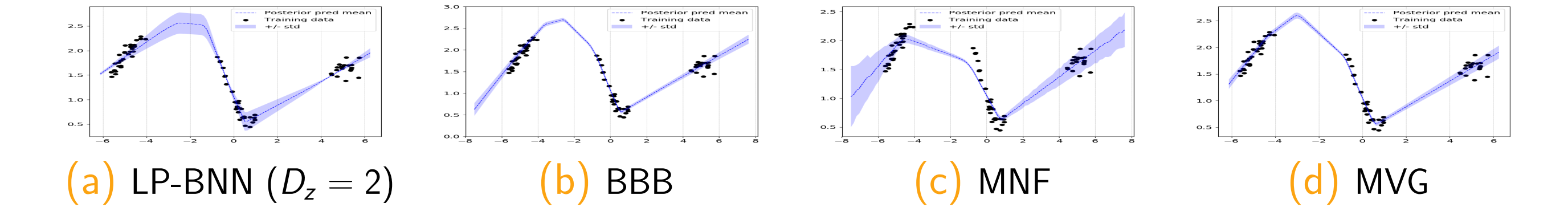


Toy data for regression is generated from a function with four modes. (left) weight space of BNN with three hidden nodes. (right) examples of functions corresponding to weights sampled from each weight cluster.

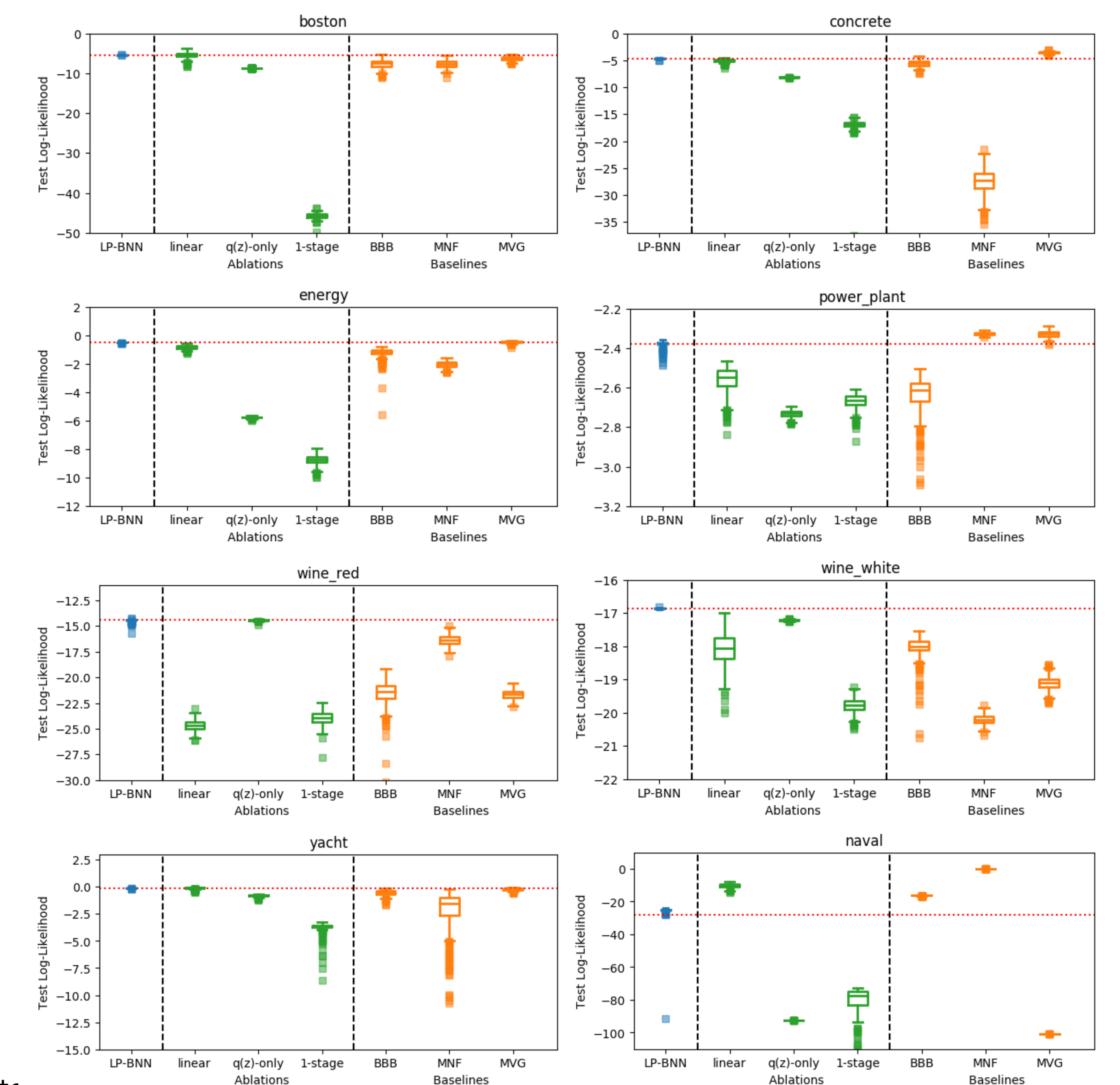


RESULTS

Baselines methods are: 1) BBB: mean field (Blundell, et.al 2015); 2) MNF: multiplicative normalizing flow (Louizos et.al, 2017); 3) MVG: multivariate Gaussian prior BNN (Louizos et.al, 2016).



Inferred predictive posterior distribution for a toy data set drawn from a NN with 1-hidden layer, 20 hidden nodes and RBF activation functions.



Test log-likelihood for UCI benchmark datasets for best dimensionality of \mathbf{z} -space. Ablations of LP-BNN are: LP-BNN, LP-BNN with linear projections (linear), LP-BNN without training the autoencoder (1-stage), LP-BNN modeling uncertainty only in \mathbf{z} ($q(\mathbf{z})$ -only).

CONCLUSION

- ▶ LP-BNN provides better uncertainty estimations
- ▶ LP-BNN achieves performs better across datasets