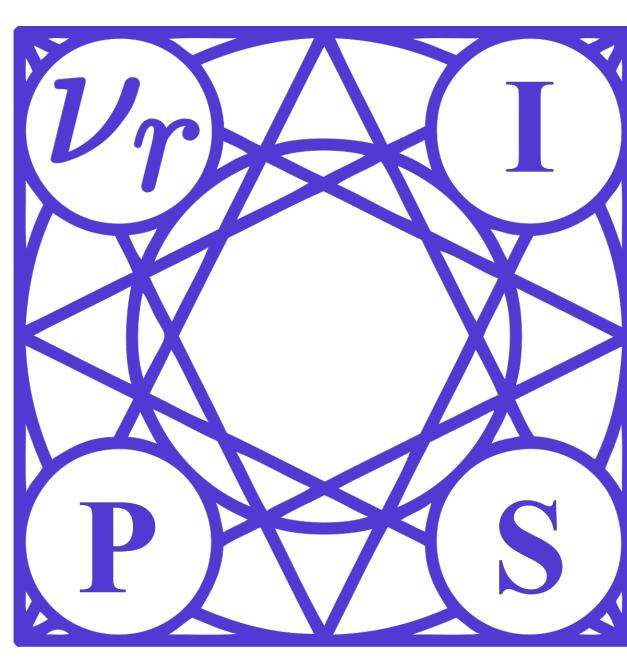


# Reliable training and estimation of variance networks

Martin Jørgensen<sup>\*,1</sup> [marjor@dtu.dk](mailto:marjor@dtu.dk), Nicki S. Detlefsen<sup>\*,1</sup> [nsde@dtu.dk](mailto:nsde@dtu.dk) and Søren Hauberg<sup>1</sup> [sohau@dtu.dk](mailto:sohau@dtu.dk)

<sup>\*</sup>Equal contribution <sup>1</sup>Technical University of Denmark



## Objective

Predictive uncertainty has not experienced the same success as mean prediction, and typically most neural network based models are overconfident. We propose methods and techniques that train *variance networks* knowing that they are different from standard regression networks.

## Methods

- ① *The locality sampler*: Estimating variance based on few observations have high variance itself. We propose a two-stage sampling scheme to replace standard mini-batching

$$\sum_{i=1}^N \left\{ -\frac{1}{2} \log(\sigma^2(\mathbf{x}_i)) - \frac{(y_i - \mu(\mathbf{x}_i))^2}{2\sigma^2(\mathbf{x}_i)} \right\} \approx \sum_{\mathbf{x}_j \in \mathcal{O}} \frac{1}{\pi_j} \left\{ -\frac{1}{2} \log(\sigma^2(\mathbf{x}_j)) - \frac{(y_j - \mu(\mathbf{x}_j))^2}{2\sigma^2(\mathbf{x}_j)} \right\}$$

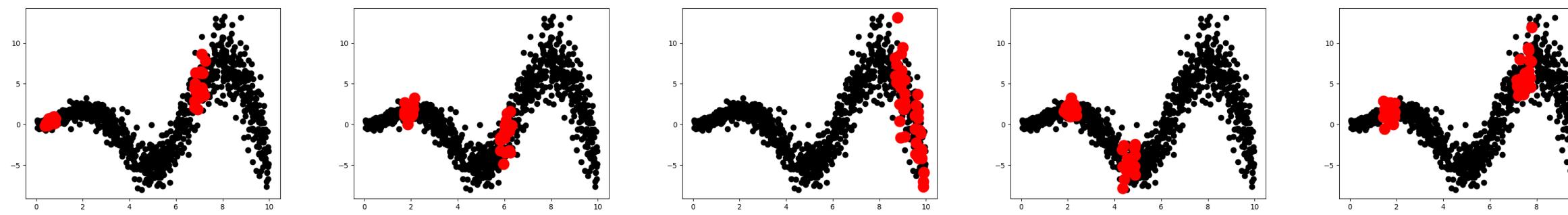


Figure: 5 batches from our proposed locality sampler

- ② *Extrapolation*: One of vanilla neural networks caveats is poor extrapolation, in the sense that predictions away from data are arbitrary. Let  $\{\mathbf{c}_i\}_{i=1}^L$  be points in  $\mathbb{R}^D$  that represent the training data, akin to inducing points in sparse GPs. Then define

$$\delta(\mathbf{x}_0) = \min_i \|\mathbf{c}_i - \mathbf{x}_0\|$$

$$\hat{\sigma}^2(\mathbf{x}_0) = (1 - \nu(\delta(\mathbf{x}_0)))\hat{\sigma}_\theta^2 + \eta\nu(\delta(\mathbf{x}_0)),$$

where  $\nu : [0, \infty) \mapsto [0, 1]$  is a surjective increasing function.

- ③ *Mean-variance split*: Ensure that the gradients are always pointing towards the MLE, by sequentially optimizing  $p_\theta(y|\theta_\mu)$  and  $p_\theta(y|\theta_{\sigma^2})$  separately.

- ④ *Student-t-likelihood* is known to be more robust in presence of small sample data – we put a Inverse-Gamma prior on  $\sigma^2$ , and instead train parameter functions  $\alpha(\mathbf{x})$  and  $\beta(\mathbf{x})$  to gain more robust estimates.

## Ablation study

We individually test each contribution, and find that they increase performance individually and in combination.

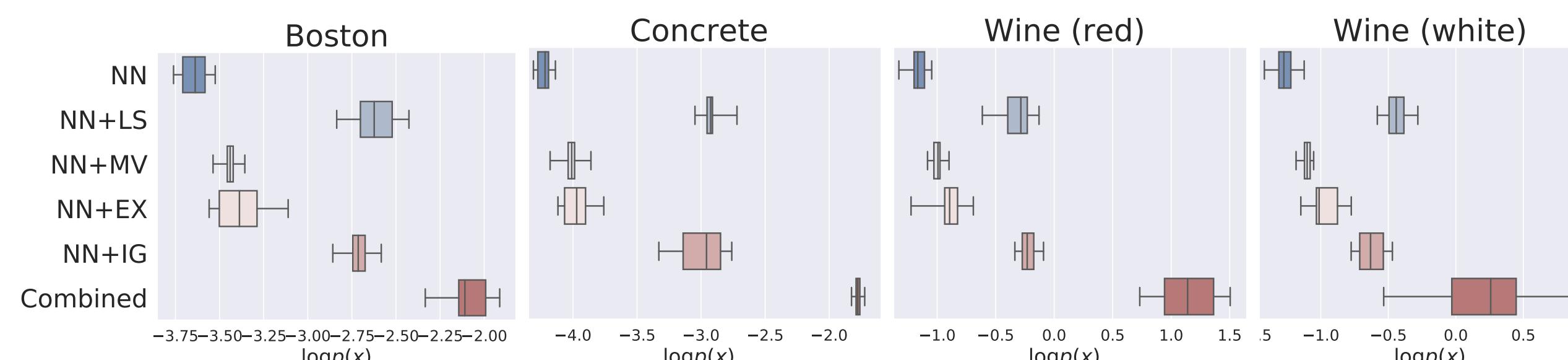


Figure: Small ablation study on four UCI benchmark datasets.

## Acknowledgements

This work was supported by a research grant (15334) from VILLUM FONDEN. This project has received funding from the European Research Council (ERC) under the European Union's Horizon 2020 research and innovation programme (grant agreement n° 757360).

## Generative Models

Variational autoencoders are among state-of-the-art models for generative models. However, these models suffer from poor uncertainty estimation when the decoder is modeled as  $p_\theta(\mathbf{x}|\mathbf{z}) = \mathcal{N}(\mathbf{x}|\mu_\theta(\mathbf{z}), \sigma_\theta^2(\mathbf{z}))$ . In the literature, we find that variance networks are avoided by using a Bernoulli distribution, modeled using a constant and often not included in generated samples. In our work we model the decoder variance using our *Combined* method.

## Toy data

We consider generative modeling of toy data. We have mapped the two half moons into four dimensions, using mappings that MLP's struggle to learn, thus successful modeling requires the variance function to be meaningful.

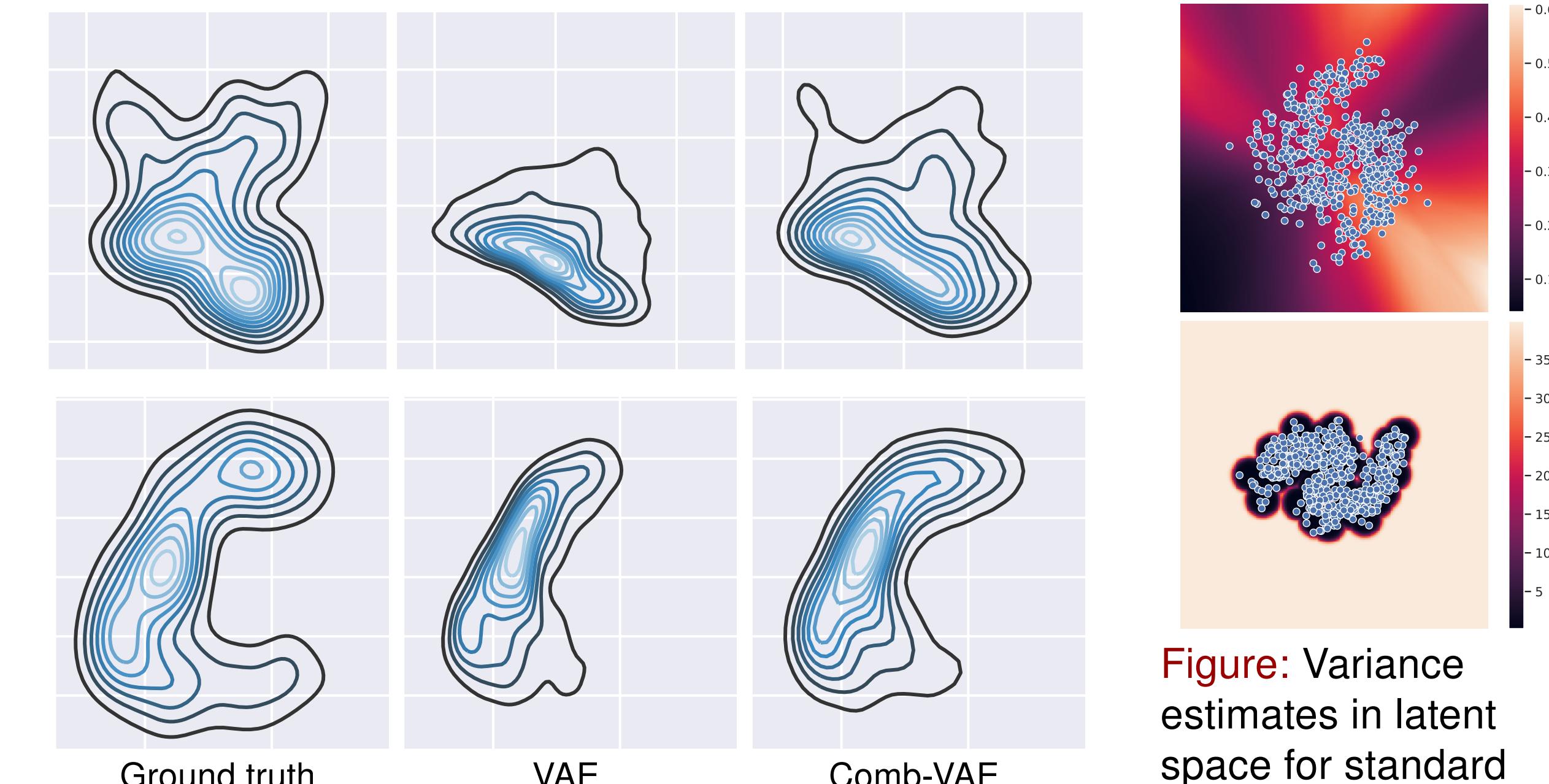


Figure: The ground truth and generated distributions. Top:  $x_1$  vs.  $x_2$ . Bottom:  $x_2$  vs.  $x_3$ .

## Image generation

We test our proposed method against standard VAE in image generation on MNIST, FashionMNIST, CIFAR and SVHN. We qualitatively get better ELBO and log-likelihood. Quantitatively we visually get a more expressive decoder variance.

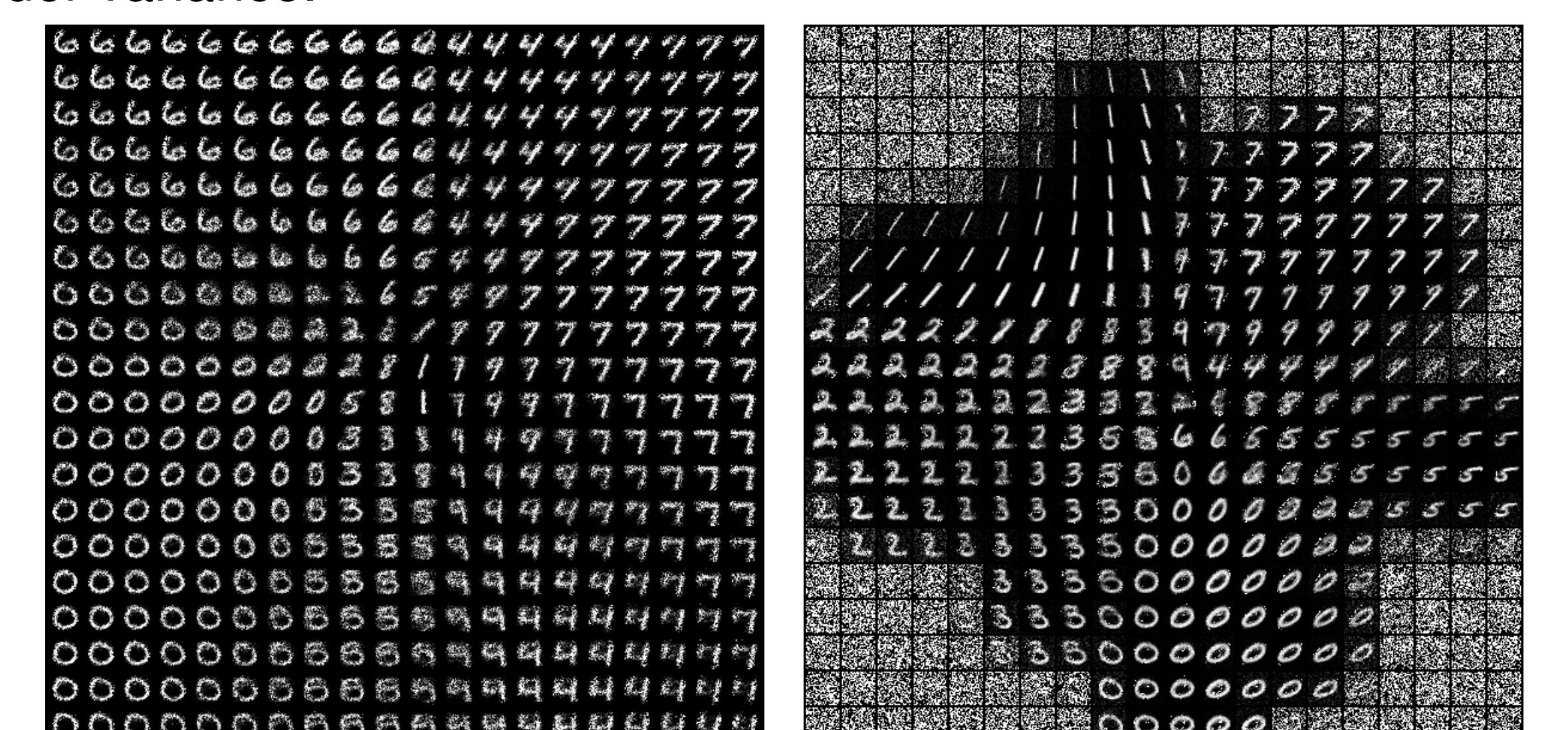


Figure: Left: Standard variance estimation are near constant and extrapolated badly. Right: Using the improved variance estimator, we are able to get more expressive variance estimated where we have data, and better extrapolation.

	MNIST	FashionMNIST	CIFAR10	SVHN
ELBO	VAE $2053.01 \pm 1.60$	$1506.31 \pm 2.71$	$1980.84 \pm 3.32$	$3696.35 \pm 2.94$
	Comb-VAE $2152.31 \pm 3.32$	$1621.29 \pm 7.23$	$2057.32 \pm 8.13$	$3701.41 \pm 5.84$
$\log p(\mathbf{x})$	VAE $1914.77 \pm 2.15$	$1481.38 \pm 3.68$	$1809.43 \pm 10.32$	$3606.28 \pm 2.75$
	Comb-VAE $2018.37 \pm 4.35$	$1567.23 \pm 4.82$	$1891.39 \pm 20.21$	$3614.39 \pm 7.91$

Table: Generative modeling of 4 image datasets.

## Regression

We consider the problem  $y = \mathbf{x} \cdot \sin(\mathbf{x}) + 0.3 \cdot \epsilon_1 + 0.3 \cdot \mathbf{x} \cdot \epsilon_2$ , with  $\epsilon_1, \epsilon_2 \sim \mathcal{N}(0, 1)$ , and compare us to Gaussian Processes (GP), Neural Networks baseline (NN), Bayesian Neural Networks (BNN), Monte Carlo Dropout (MC-Dropout) and Deep Ensembles (Ens-NN).

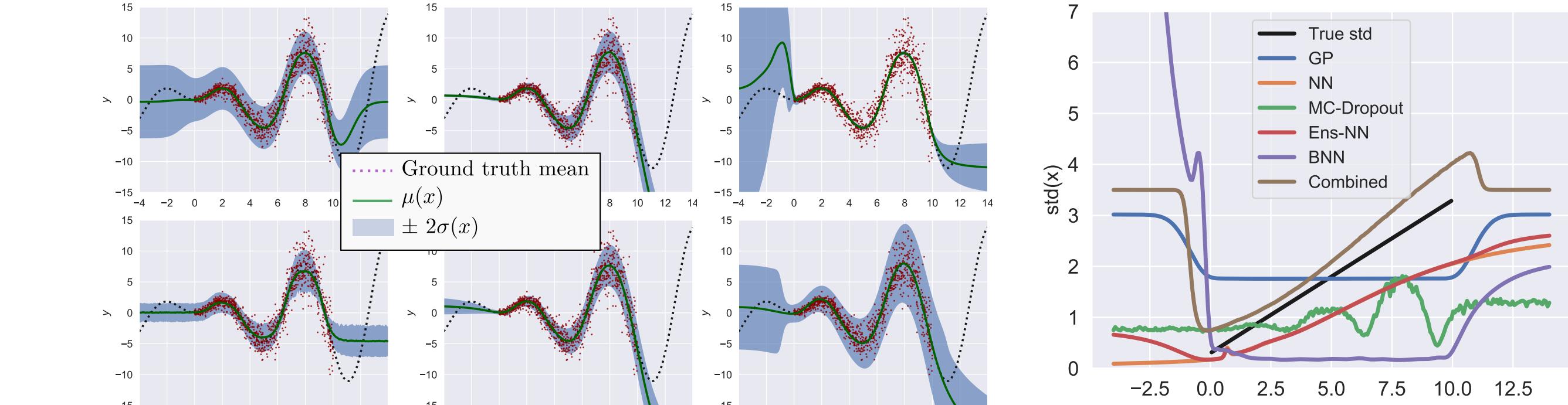


Figure: From top left to bottom right: GP, NN, BNN, MC-Dropout, Ens-NN, Combined.

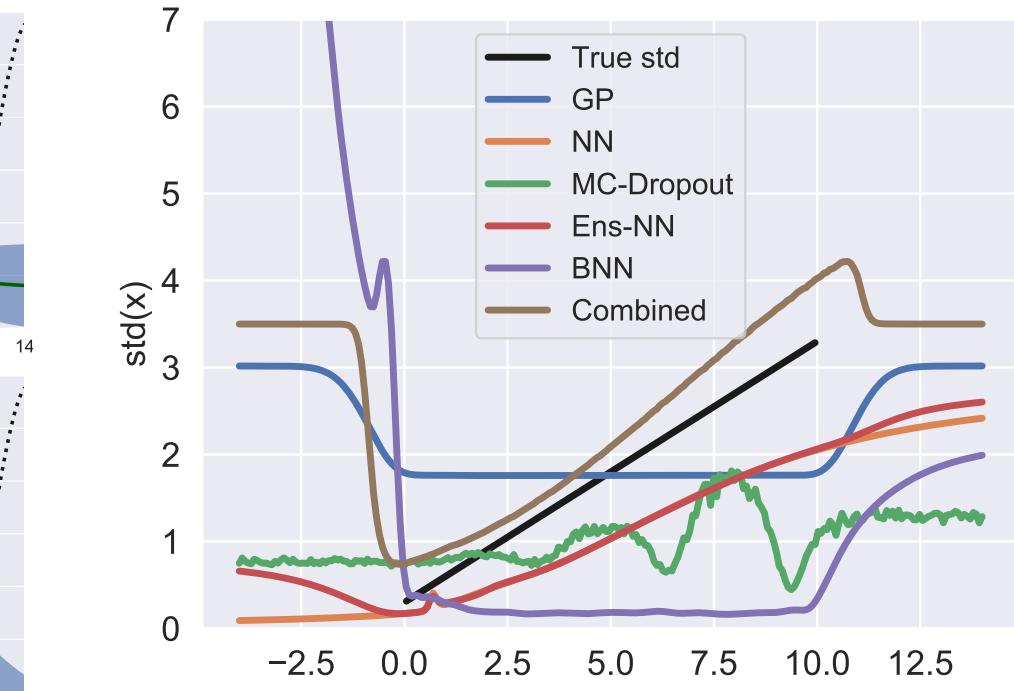


Figure: Standard deviation estimates as a function of  $x$ .

## UCI Benchmark

We outperform current methods, on 9 out of 13 dataset measured by test set log likelihood

	$N$	$D$	GP	SGP	NN	BNN	MC-Dropout	Ens-NN	Combined
Boston	506	13	$-1.76 \pm 0.3$	$-1.85 \pm 0.25$	$-3.64 \pm 0.09$	$-2.59 \pm 0.11$	$-2.51 \pm 0.31$	$-2.45 \pm 0.25$	$-2.09 \pm 0.09$
Carbon	10721	7	-	$3.74 \pm 0.53$	$2.03 \pm 0.14$	$-1.1 \pm 1.76$	$-1.08 \pm 0.05$	$-0.44 \pm 7.28$	$4.35 \pm 0.16$
Concrete	1030	8	$-2.13 \pm 0.14$	$-2.29 \pm 0.12$	$-2.3 \pm 0.07$	$-3.11 \pm 0.12$	-	$-3.06 \pm 0.32$	$-1.78 \pm 0.04$
Energy	768	8	$-1.85 \pm 0.34$	$-2.22 \pm 0.15$	$-3.78 \pm 0.04$	$-2.07 \pm 0.08$	$-2.01 \pm 0.11$	<b><math>-1.48 \pm 0.31</math></b>	$-1.68 \pm 0.13$
Kin8nm	8192	8	-	$2.01 \pm 0.02$	$0.08 \pm 0.02$	$0.95 \pm 0.08$	$0.95 \pm 0.15$	$1.18 \pm 0.03$	$2.49 \pm 0.07$
Naval	11934	16	-	-	$3.47 \pm 0.21$	$3.71 \pm 0.05$	$3.80 \pm 0.09$	$5.55 \pm 0.05$	$7.27 \pm 0.13$
Power plant	9568	4	-	$-1.9 \pm 0.03$	$-4.26 \pm 0.14$	$-2.89 \pm 0.01$	$-2.89 \pm 0.14$	$-2.77 \pm 0.04$	$-1.19 \pm 0.03$
Protein	45730	9	-	-	$-2.95 \pm 0.09$	$-2.91 \pm 0.00$	$-2.93 \pm 0.14$	<b><math>-2.80 \pm 0.02</math></b>	$-2.83 \pm 0.05$
Superconduct	21263	81	-	$-4.07 \pm 0.01$	$-4.92 \pm 0.10$	$-3.06 \pm 0.14$	$-2.91 \pm 0.19$	$-3.01 \pm 0.05$	$-2.43 \pm 0.05$
Wine (red)	1599	11	$0.96 \pm 0.18$	$-0.08 \pm 0.01$	$-0.11 \pm 0.098$	$-0.01 \pm 0.04$	$-0.94 \pm 0.01$	$-0.93 \pm 0.09$	$1.21 \pm 0.23$
Wine (white)	4898	11	-	$-0.14 \pm 0.05$	$-1.29 \pm 0.09$	$-0.41 \pm 0.17$	$-1.26 \pm 0.01$	$-0.99 \pm 0.06$	$0.40 \pm 0.42$
Yacht	308	7	$0.16 \pm 1.22$	$-0.38 \pm 0.32$	$-4.12 \pm 0.17$	$-1.65 \pm 0.05$	$-1.55 \pm 0.12$	$-1.18 \pm 0.21$	$-0.07 \pm 0.05$
Year	51534590	-	-	-	$-5.21 \pm 0.87$	$-3.97 \pm 0.34$	$-3.78 \pm 0.01$	$-3.42 \pm 0.02$	$-3.01 \pm 0.14$

Table: Dataset characteristics and tests set log-likelihoods for the different methods. A “-” indicates the model was infeasible to train. Bold highlights the best results.

## Active Learning

We construct a active learning setting based on UCI datasets. On some dataset, we see observe faster learning than other methods. On other datasets, we are on par with existing methods.

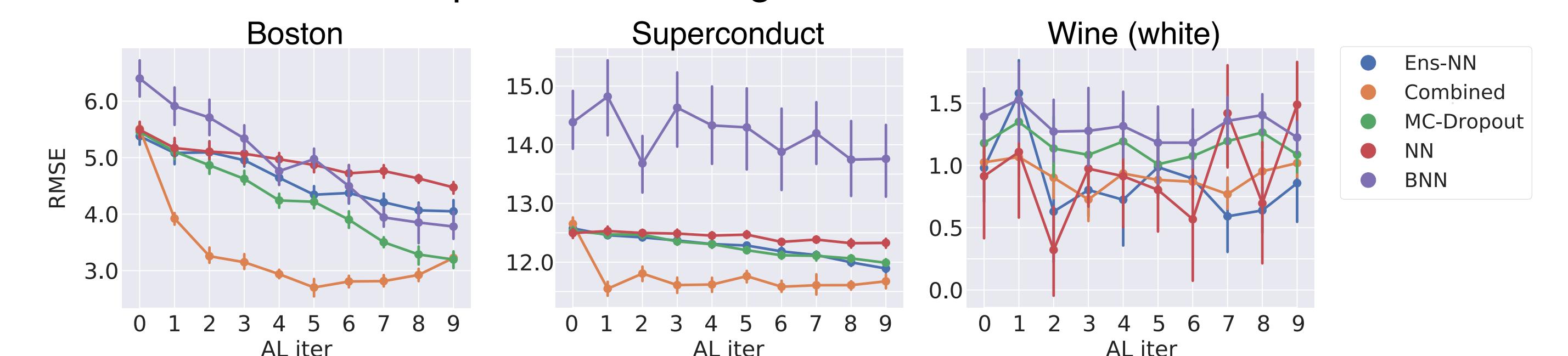


Figure: Average test set RMSE and standard errors in active learning.

## Conclusions

- We propose methodologies for better variance estimation in neural nets, that can be used in combination or individually.
- We clearly outperform established ensemble methods in regression modelling and active learning.
- In generative modeling, we better capture the data distribution and get more interpretable variance estimates than baseline methods.