

This paper describes objective technical results and analysis. Any subjective views or opinions that might be expressed in the paper do not necessarily represent the views of the U.S. Department of Energy or the United States Government.

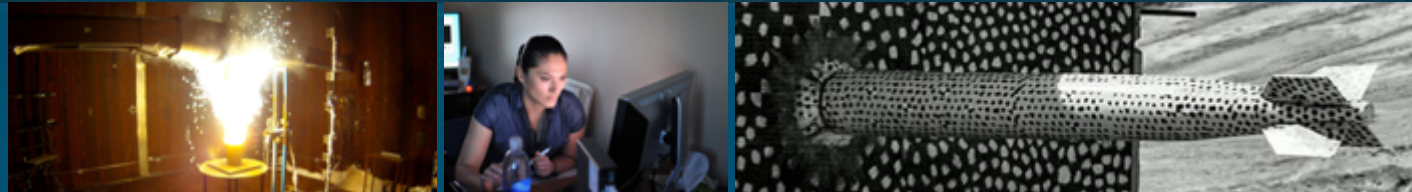


Sandia
National
Laboratories

SAND2020-12732C

Estimating Predictive Uncertainty in Scientific Machine Learning

A Library of Methods and Test Problems



Ahmad A. Rushdi, Aubrey C. Eckert, Gabriel Huerta,
Laura P. Swiler, Brian M. Adams

Sandia National Laboratories

Workshop
Uncertainty Management
and Machine Learning
in Engineering
Applications
November 16-17, 2020



Sandia National Laboratories is a multimission laboratory managed and operated by National Technology & Engineering Solutions of Sandia, LLC, a wholly owned subsidiary of Honeywell International Inc., for the U.S. Department of Energy's National Nuclear Security Administration under contract DE-NA0003525.



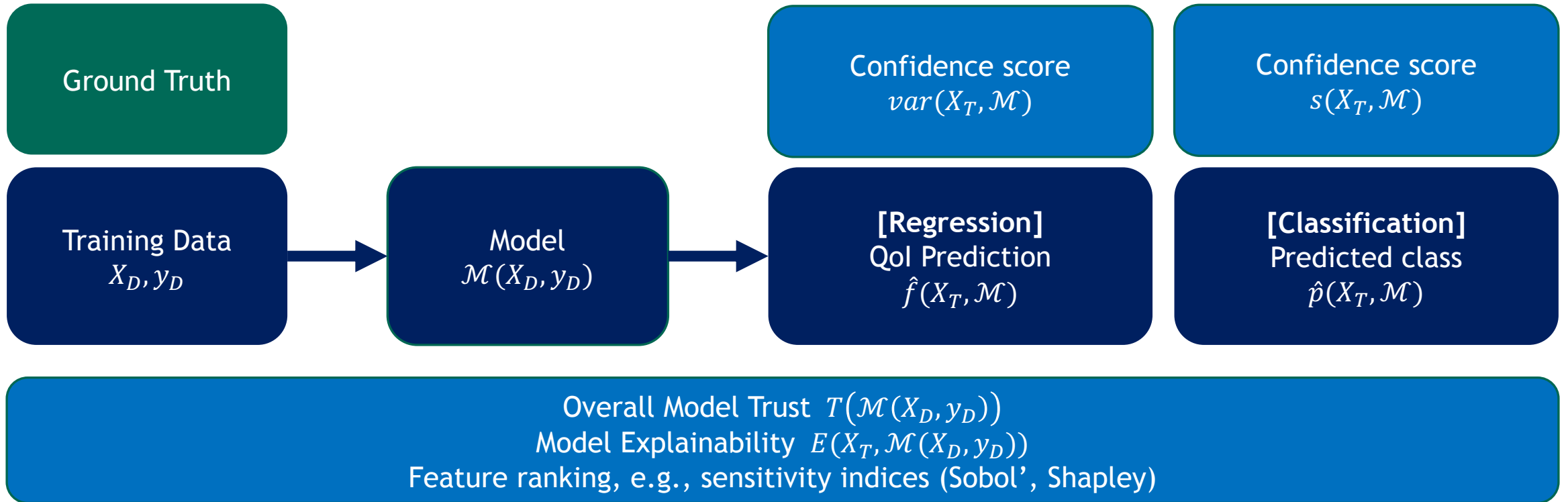
1. Problem Statement
2. Why Predictive Uncertainty?
3. Analytical Library of Test Functions
4. Methods
5. Numerical Experiments
 - Regression
 - Classification
6. Concluding Remarks



Neural network models have attracted a lot of research attention in Scientific Machine Learning (SciML) problems. However, they tend to be overconfident when reporting typical point-estimate predictions in classification and regression problems. This could be very harmful when dealing with costly numerical simulations or high-stakes decisions in national security applications.

In this work, we assess uncertainty quantification techniques for neural network models. To understand their variability, we rely on different sources of randomness associated with training samples, weight initialization, dropout methods, and ensemble formations. Motivated by typical SciML situations, we assume a limited sample budget, noisy training data, and suggest approaches for reporting and possibly reducing uncertainty.

If an NN trained on some dataset is evaluated on a totally different dataset, it should be able to report higher predictive uncertainty. Inputs from a different dataset would be far away from the training data.



Supervised Machine Learning: extract *models* from *data* and use them to make predictions.

Given data: $\mathcal{D} = \{(\mathbf{x}_1, y_1), \dots, (\mathbf{x}_N, y_N)\} \subset \mathcal{X} \times \mathcal{Y}$
 minimize mean squared error of: $\frac{1}{N} \sum_i \|y_i - (\mathbf{x}_i \mathbf{W} + \mathbf{b})\|^2$

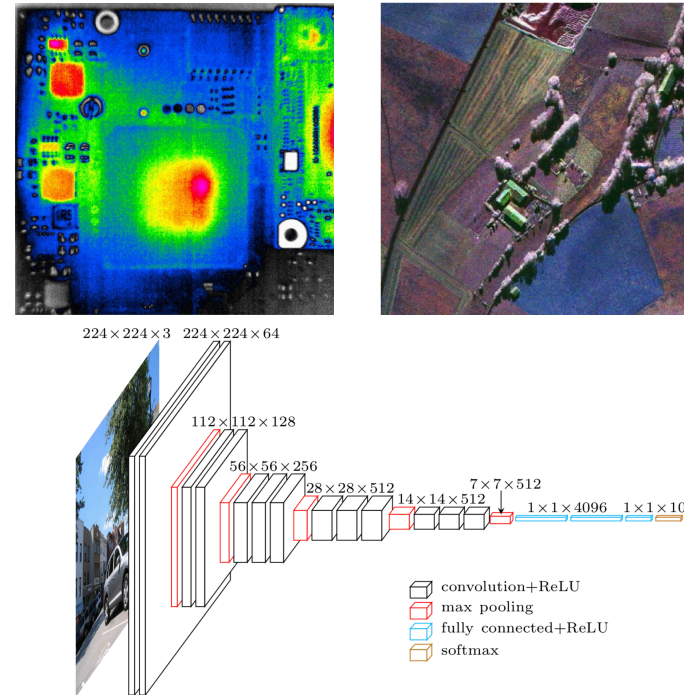


Deep learning “black box” models are popular, yet are difficult to interpret and understand.

UQ today underpins many decision processes in nuclear security, our risk management and associated investments, which can be at the scale of billions of dollars. Predictions without UQ are neither predictions nor actionable. The data-rich world of ML, especially the powerful deep learning (DL) models, poses parallel challenges.

To develop consequential decision support from ‘learned’ models built on complex datasets, there is an important need to co-develop UQ for this domain.

Feature engineering in deep learning: rethinking architectures



- Begoli E, Tanmoy B, and Dimitri K "The need for uncertainty quantification in machine-assisted medical decision making." Nature Machine Intelligence, 2019
- Mehta, Pankaj, et al. "A high-bias, low-variance introduction to machine learning for physicists." Physics reports 810 (2019): 1-124

6 Predictive Uncertainty



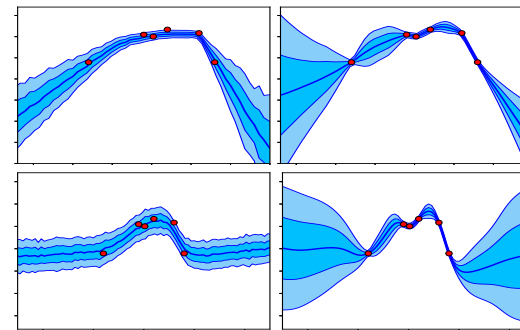
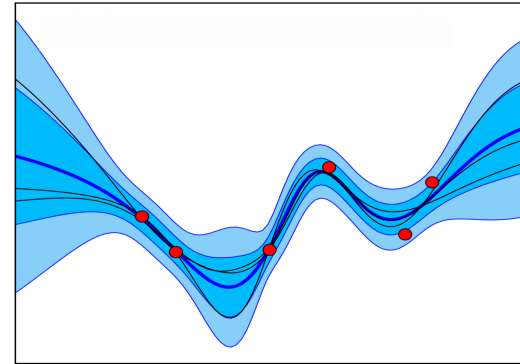
We represent DL predictions as **distributions** (μ, σ^2) rather than point estimates \hat{p} .

Scientific (vs. data rich) Machine Learning

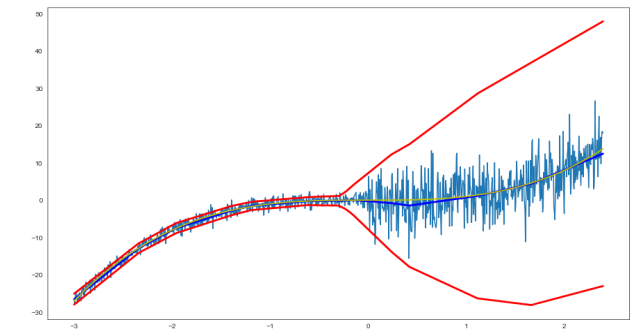
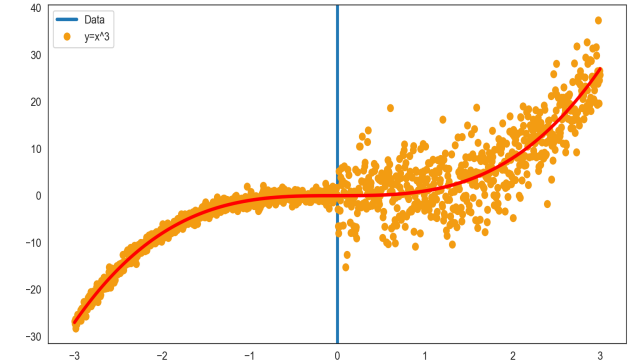
1. Training sample **budget is limited**, causing Out of Distribution (OoD) issues.
2. Different **types of noise**, imposed onto the input data and propagate that noise into an uncertainty metric in the resulting prediction.

Active Learning / Adaptive Sampling

- Model chooses which unlabeled data are most informative
- An uncertainty estimator (e.g., variance) suggests points of highest value for model accuracy improvement



Very few training samples



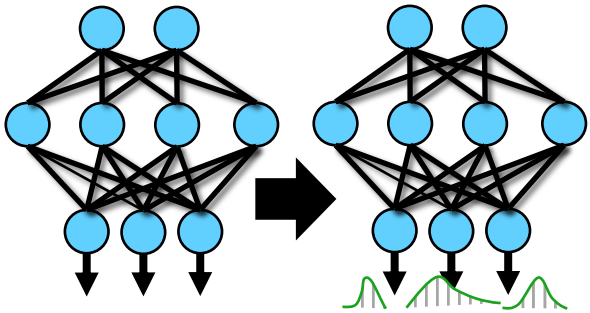
Noisy training data

Quality of prediction measures

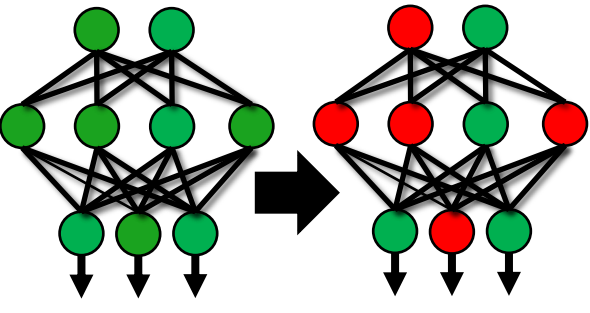
Our approach is to assign a numerical score to a prediction $p_{\theta}(y|\mathbf{x})$, rewarding *agreed upon predictions* over worse.



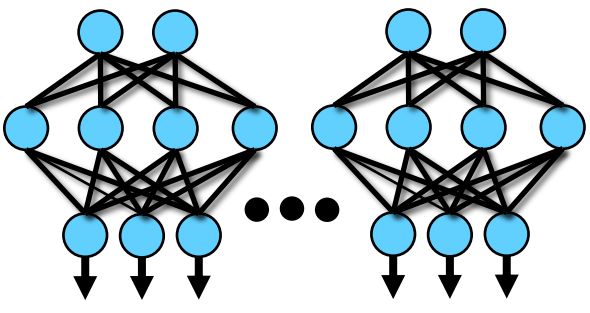
Bayesian NNs



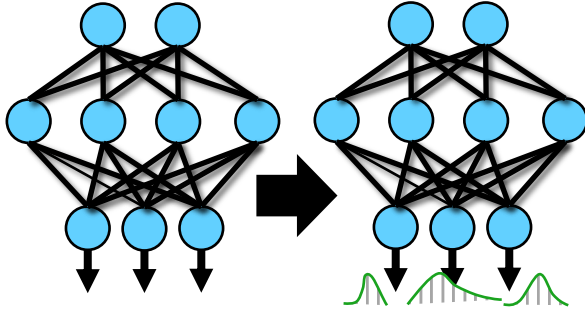
MC Dropout



Deep Ensembles



Bayesian NNs



Bayesian Neural Networks (BNNs): place a prior distribution over the network weights and use data to learn a posterior distribution.

$$\begin{aligned} \mathbf{w}^{\text{MLE}} &= \arg \max_{\mathbf{w}} \log p(\mathcal{D}|\mathbf{w}) \\ &= \arg \max_{\mathbf{w}} \sum_i \log p(\mathbf{Y}_i|\mathbf{X}_i, \mathbf{w}) \\ \mathbf{w}^{\text{MAP}} &= \arg \max_{\mathbf{w}} \log p(\mathbf{w}|\mathcal{D}) \\ &= \arg \max_{\mathbf{w}} \log p(\mathcal{D}|\mathbf{w}) + \log p(\mathbf{w}) \end{aligned}$$

Variational Inference (VI) approximation:

$$\theta^* = \arg \min KL(q_{\theta}(\mathbf{w}|\mathcal{D})||p(\mathbf{w}|\mathcal{D}))$$

Exact Bayesian inference of posteriors:
computationally intractable.

Approximations:

- Markov chain Monte Carlo (MCMC)
- Laplace approximation
- Hamiltonian methods
- Variational Bayesian methods

Quality assessment:

- Choice of priors
- Degree of approximation

Hurdles:

Computationally slower to train. Harder to implement.

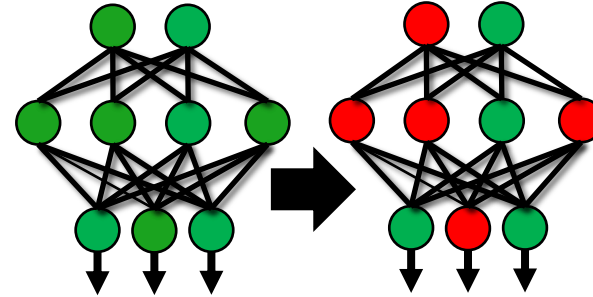


Approximate Bayesian Interpretation

Popular in practice

Simplicity of implementation at test time.

MC Dropout



Neurons are randomly dropped in each iteration, with some probability (p), often fixed at empirical values (e.g., 0.3).

$$\mathbb{E}_q(y^* | x^*)(y^*) \approx \frac{1}{T} \sum_{t=1}^T \hat{y}^*(x^*, W_1^t, \dots, W_L^t)$$

Average of T stochastic forward passes through the network.

Dropout may be interpreted as an ensemble:

Predictions are averaged over an ensemble of NNs, when dropout rates are not tuned based on training.

Y. Gal and Z. Ghahramani. Dropout as a Bayesian approximation: Representing model uncertainty in deep learning. In ICML, 2016.

Ensembles of NNs

have been successfully used to boost predictive performance in traditional ML problems (e.g. classification accuracy in ImageNet).

Variance

MSE captures only predictive mean, but not variance.

Adversarial training

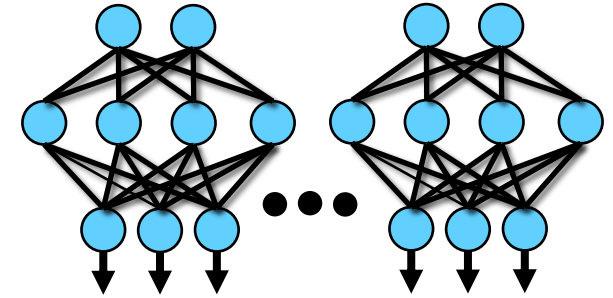
Typically: improves robustness to adversarial examples

Fast Gradient Sign Method (FGSM):

$$\mathbf{x}' = \mathbf{x} + \epsilon \text{sign}(\nabla_{\mathbf{x}} \ell(\theta, \mathbf{x}, y))$$

Here: smooth predictive distributions.

Deep Ensembles



Estimate 2 outputs: mean and variance

Modified loss function + adversarial training

$$-\log p_{\theta}(y_n | \mathbf{x}_n) = \frac{\log \sigma_{\theta}^2(\mathbf{x})}{2} + \frac{(y - \mu_{\theta}(\mathbf{x}))^2}{2\sigma_{\theta}^2(\mathbf{x})}$$

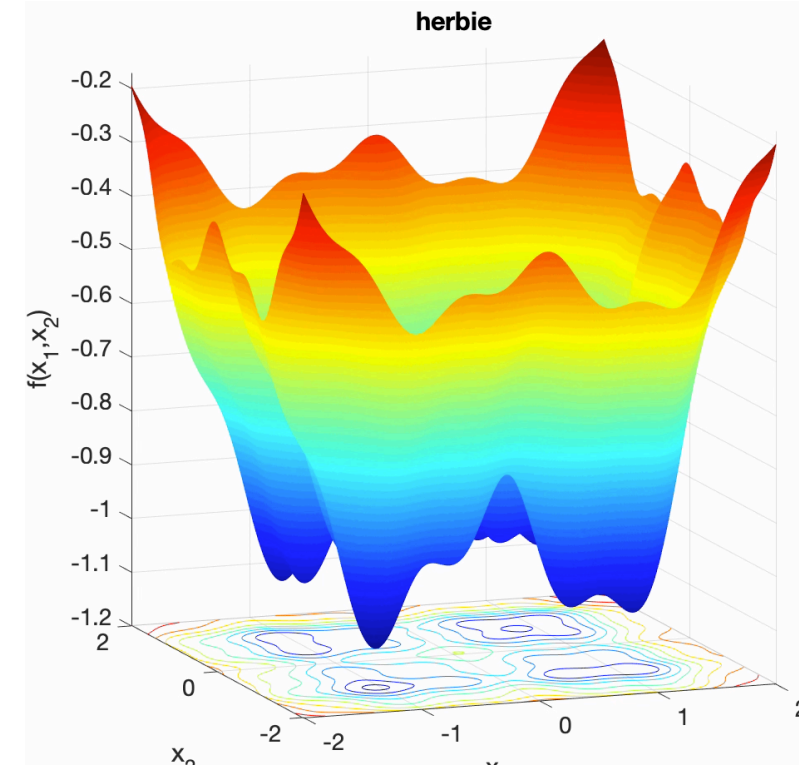
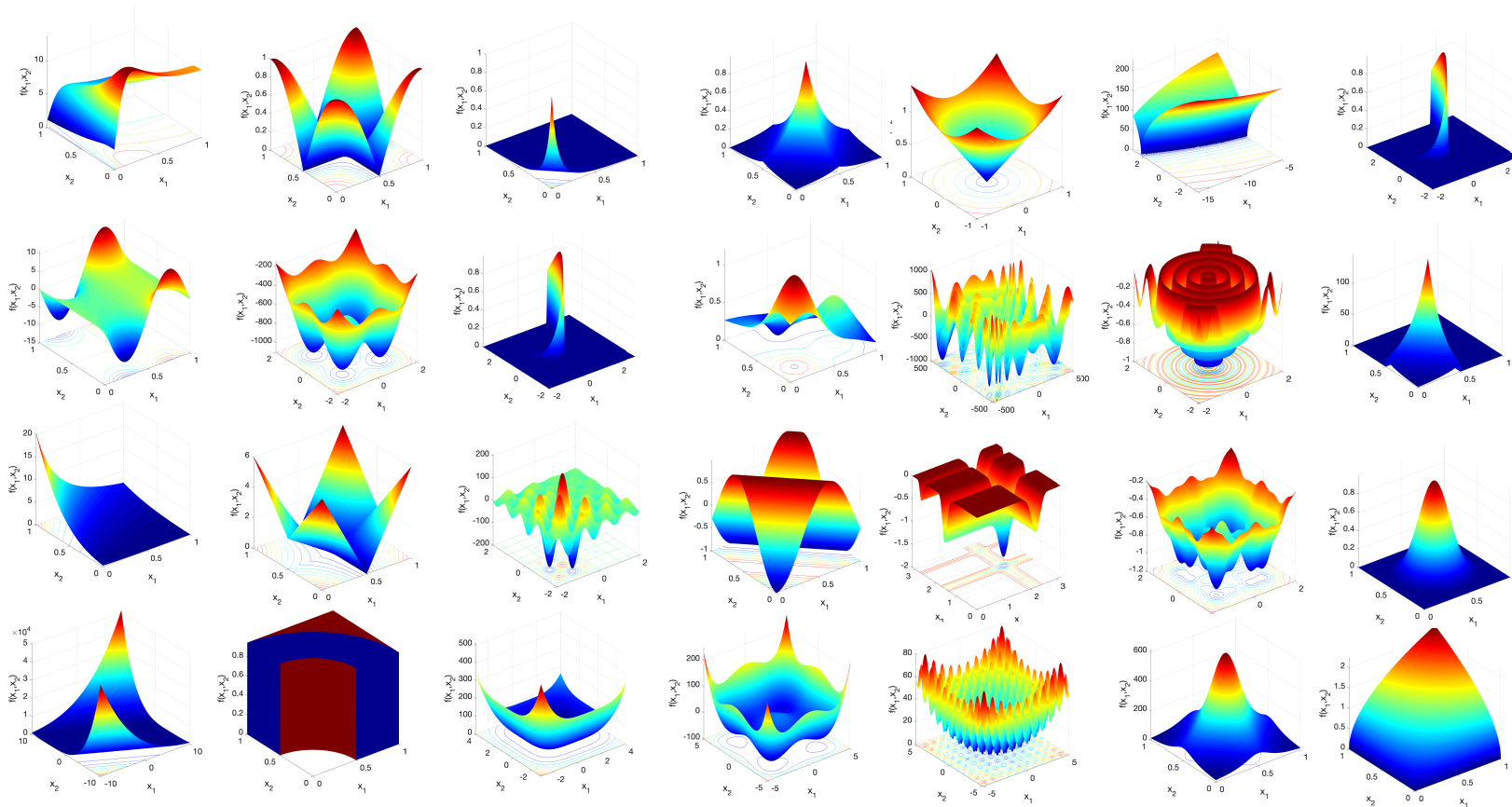
Decomposed uncertainty

$$\mu_e = \frac{1}{M} \sum_{i=1}^M \mu_i \text{ and}$$

$$\sigma_e^2 = \underbrace{\frac{1}{M} \sum_{i=1}^M \sigma_i^2}_{\text{Aleatoric}} + \underbrace{\left[\frac{1}{M} \sum_{i=1}^M \mu_i^2 - \mu_e^2 \right]}_{\text{Epistemic}}$$



A library of **high-dimensional analytical test functions**, frequently used for QoI experimental studies (optimization, numerical integration, uncertainty quantification, and multi-fidelity analysis).



- Dakota www.Dakota.sandia.gov.
- Virtual Library of Simulation Experiments <https://www.sfu.ca/~ssurjano/index.html>

Numerical Experiments - Regression



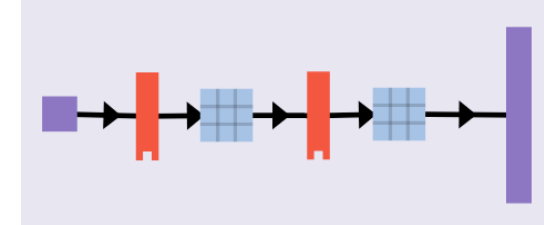
Data:

function samples + $\begin{cases} \mathcal{AWGN}(0, \sigma_1) \\ \mathcal{AWGN}(0, \sigma_2) \end{cases}$

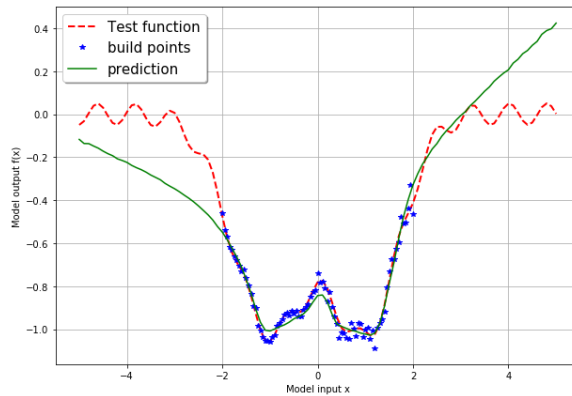
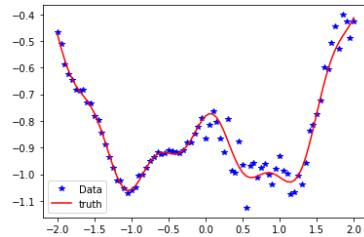
Model:

```
input = Input(shape=(1,))
x = Dense(512, activation="relu")(input)
x = Dropout(0.5)(x, training=True)
x = Dense(512, activation="relu")(x)
x = Dropout(0.5)(x, training=True)
output = Dense(1)(x)
```

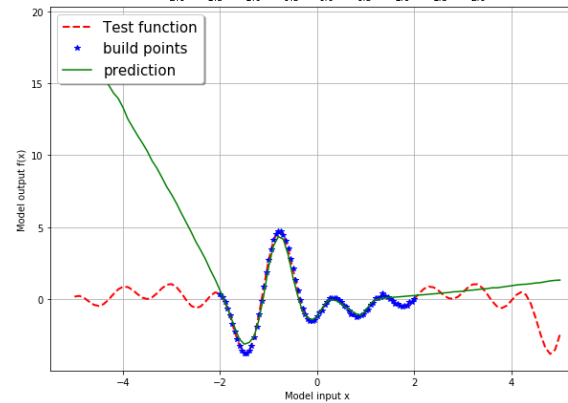
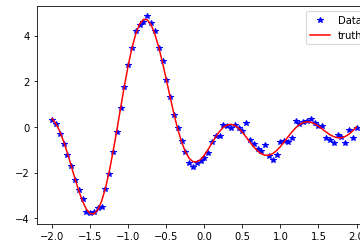
Total params: 264,193



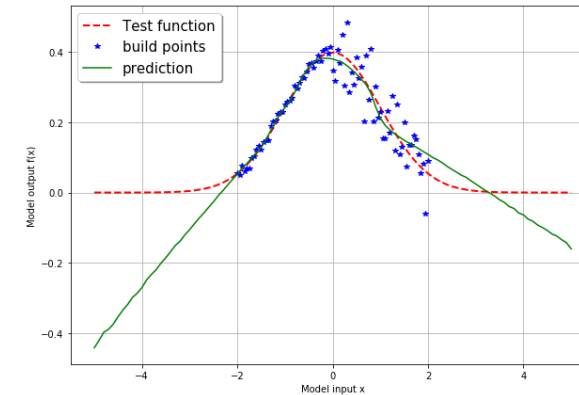
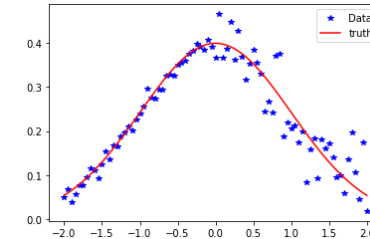
Herbie



Shubert



Gaussian

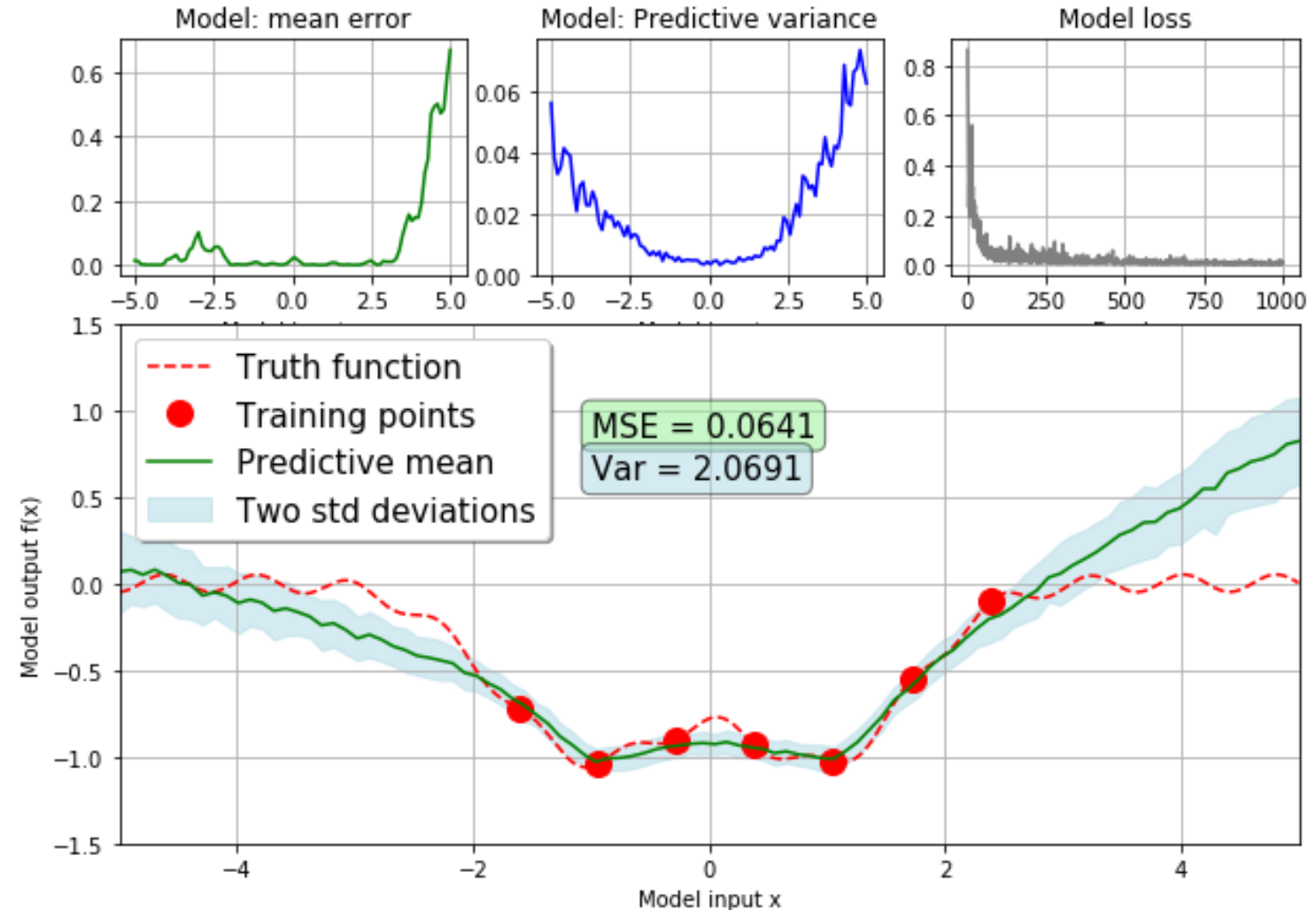


A model with 200K trainable parameters still needs an assessment of confidence/uncertainty.



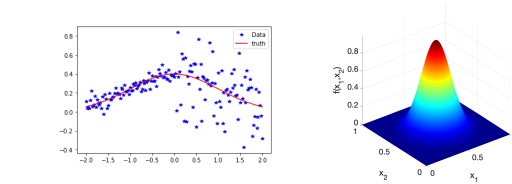
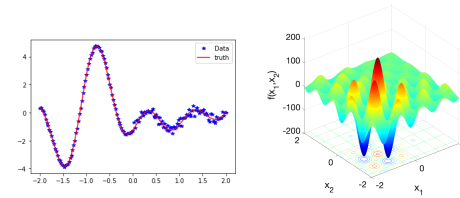
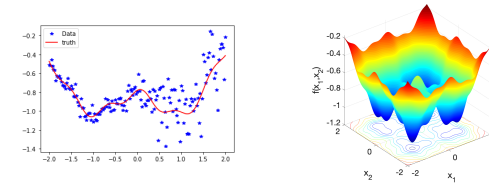
Global Sensitivity Analysis for Model Explainability

- Uncertainty quantification at a new **test point**.
- Overall model uncertainty/confidence within a **test domain**.
- Guidance of **adaptive sampling** towards points/regions of high estimated variance.

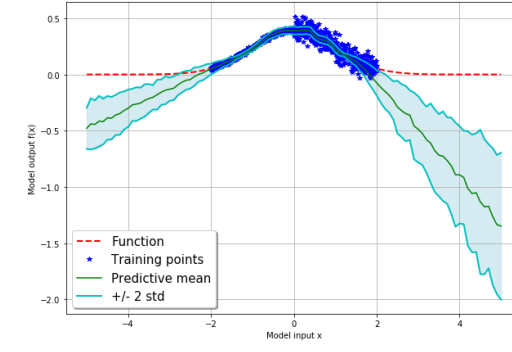
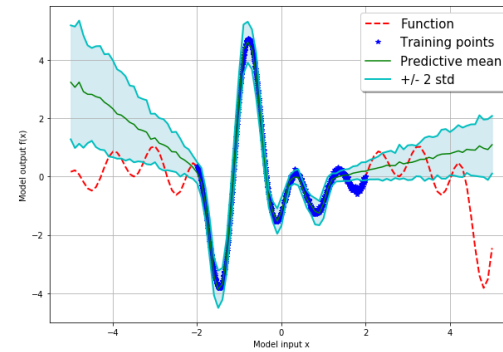
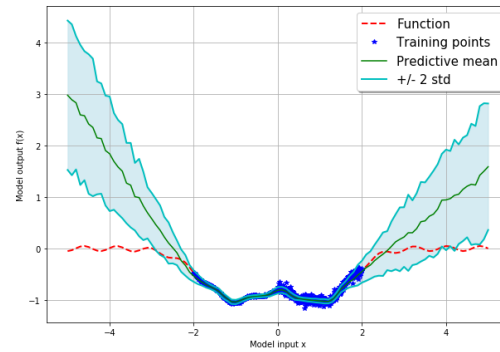




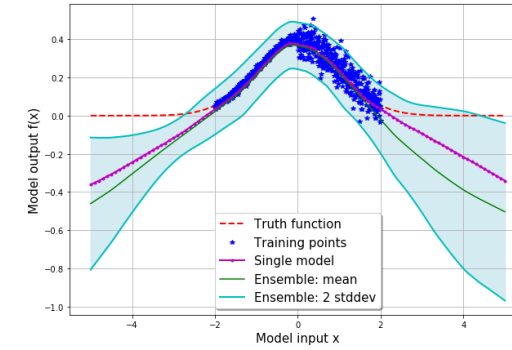
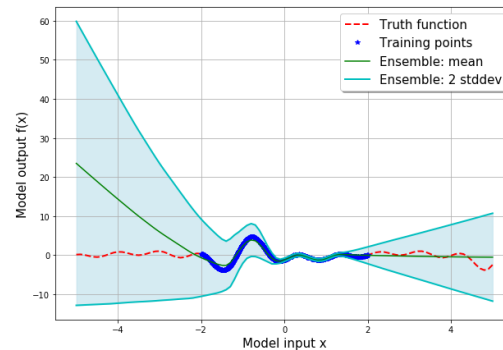
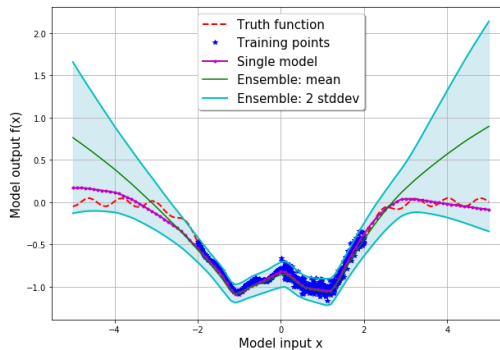
Analytical
Functions



Dropout



Deep Ensembles



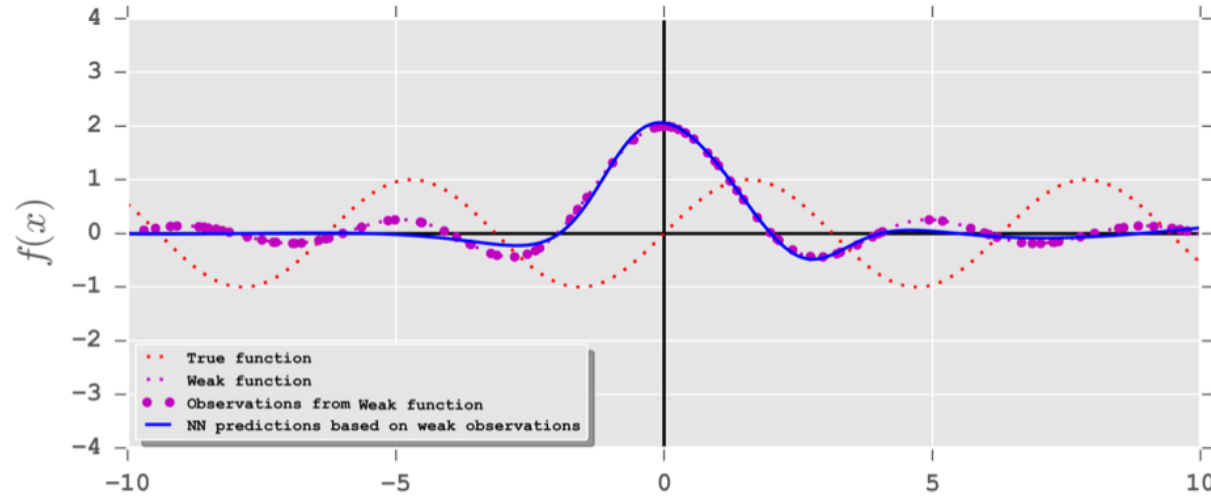
Methods vary in terms of complexity, accuracy, scalability, and computational cost.

Weighted Fidelity Learning

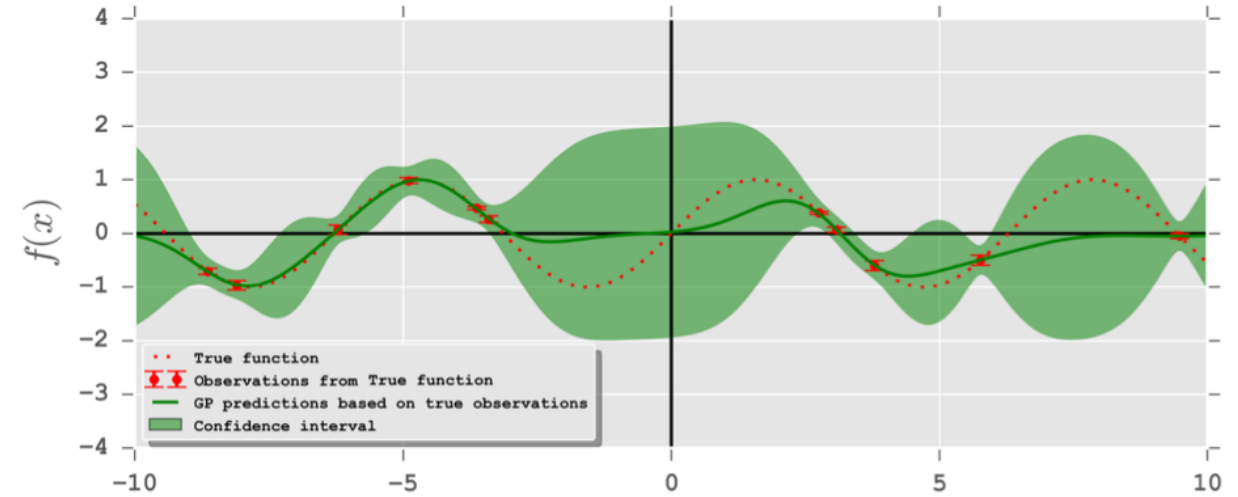


True/HF: $y = \sin(x)$
 Weak/LF: $y = 2 \operatorname{sinc}(x)$

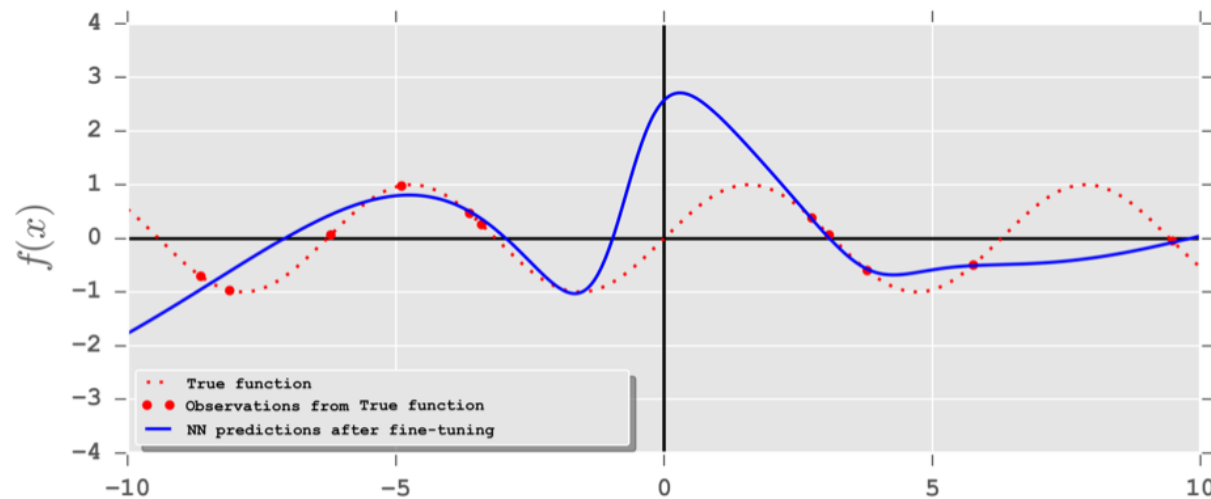
Training LF-NN on 100 examples from the weak function.



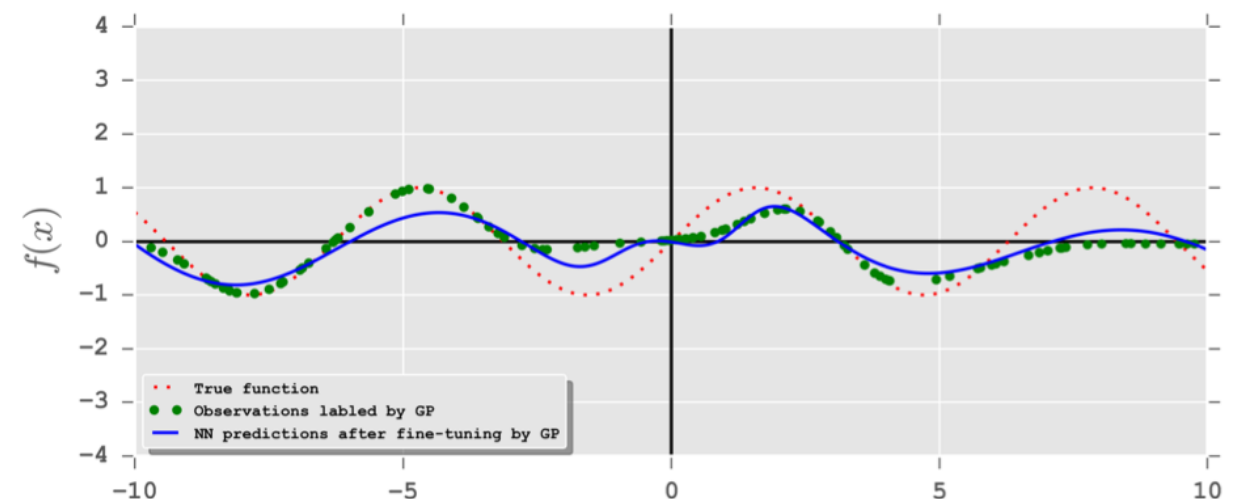
Fitting HF-NN based on 10 observations from the true function.



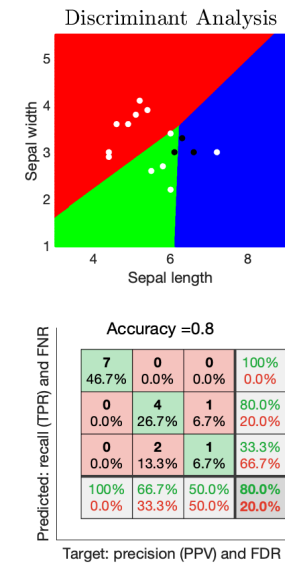
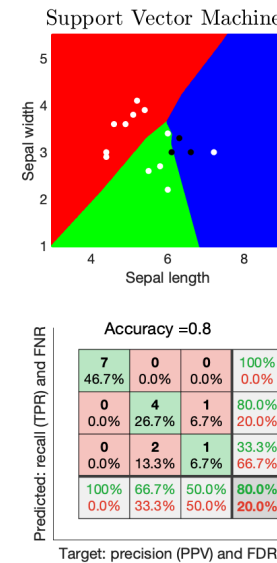
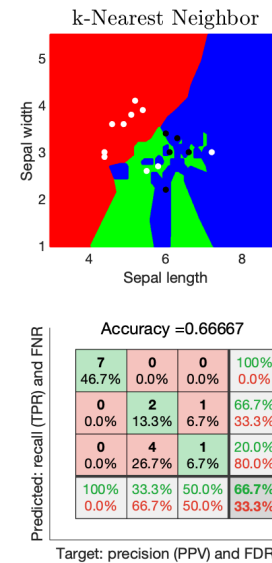
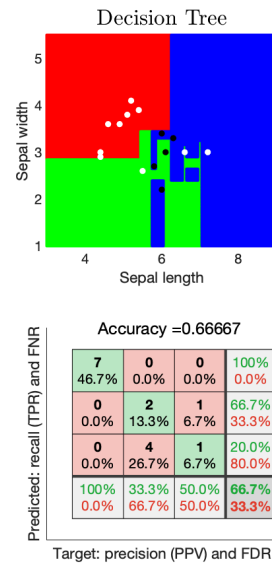
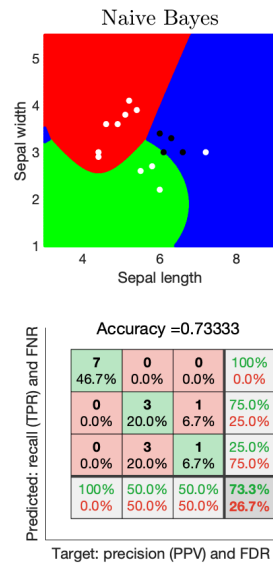
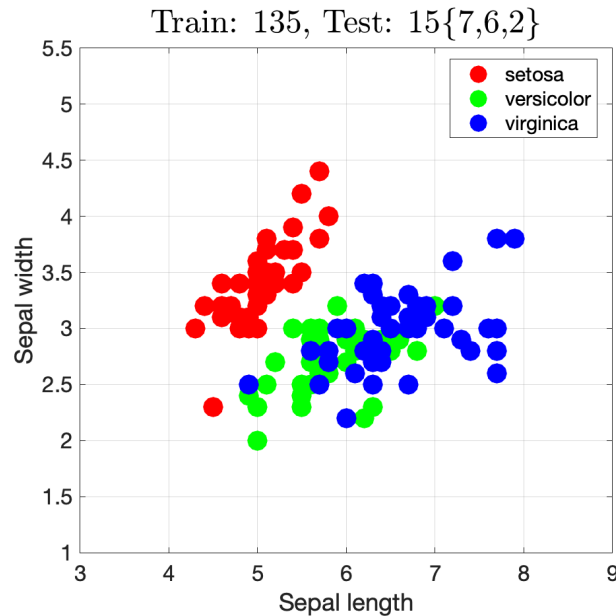
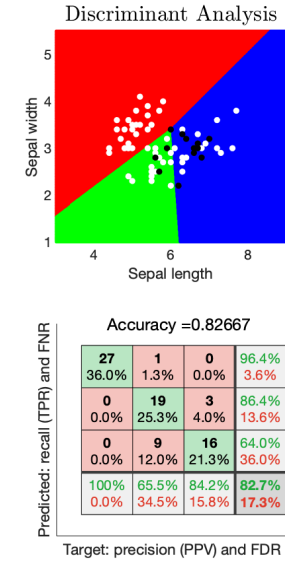
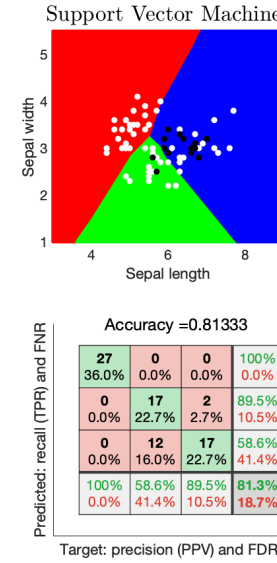
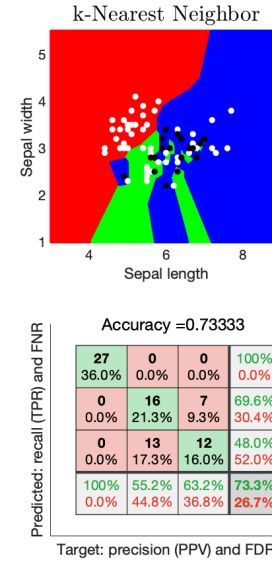
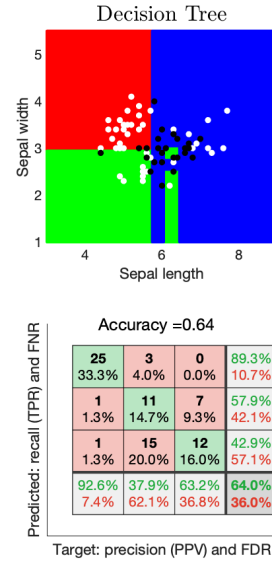
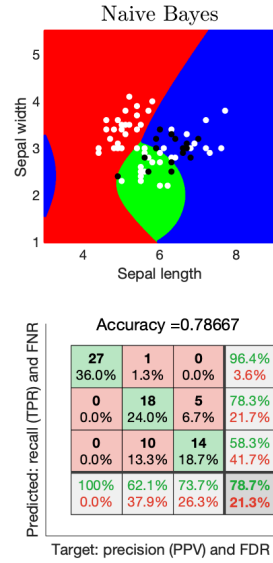
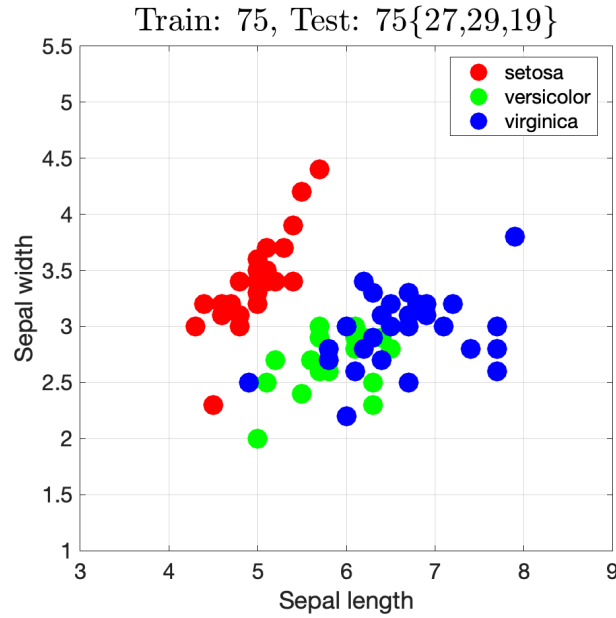
Fine-tuning LF-NN based on observations from the true function.



Fine-tuning LF-NN based on label/confidence from HF-NN.



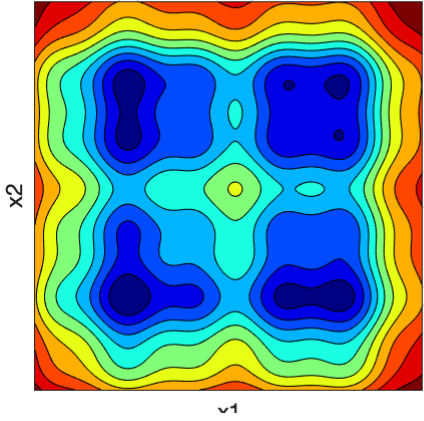
Numerical Experiments – IRIS Classification



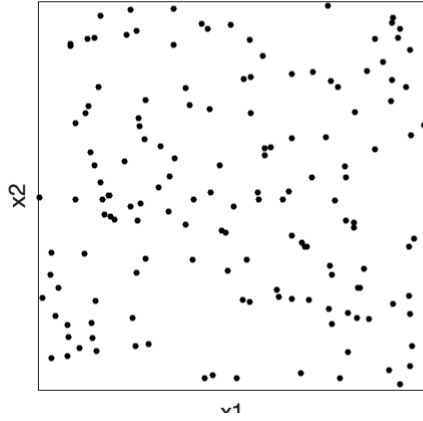
Numerical Experiments – Classification on x



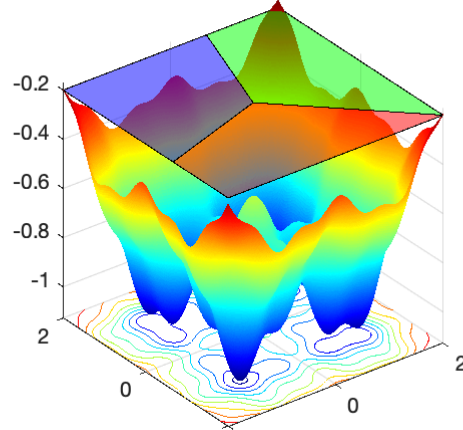
Underlying Behavior



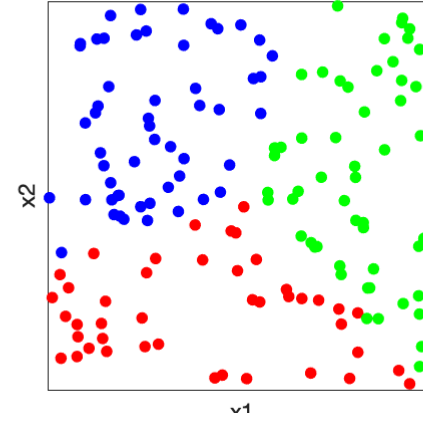
Training Samples (MC)



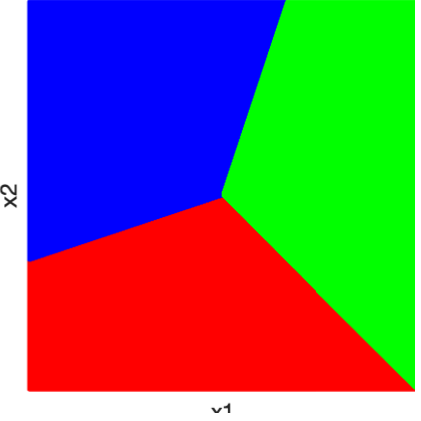
Labeling



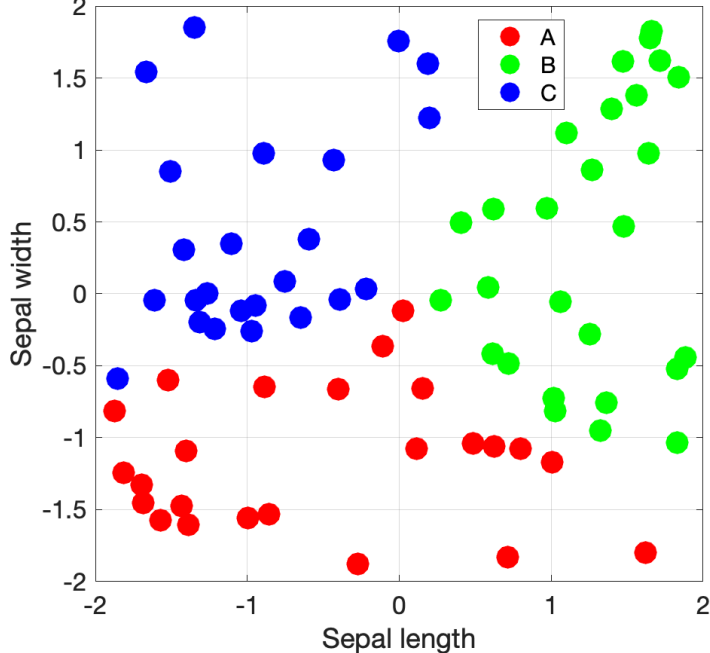
Ground Truth Labels



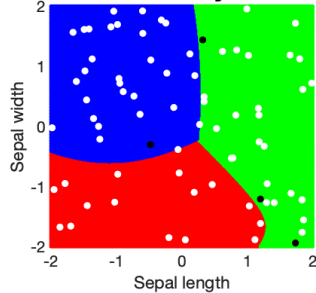
Ideal Classifier



Train: 75, Test: 75{19,27,29}



Naive Bayes

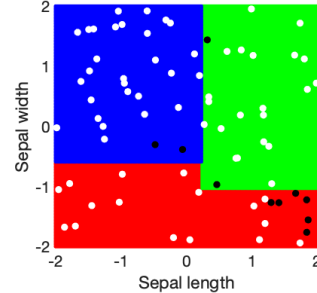


Accuracy = 0.94667

16	0	0	100%
21.3%	0.0%	0.0%	0.0%
2	27	1	90.0%
2.7%	36.0%	1.3%	10.0%
1	0	28	96.6%
1.3%	0.0%	37.3%	3.4%
84.2%	100%	96.6%	94.7%
15.8%	0.0%	3.4%	5.3%

Target: precision (PPV) and FDR

Decision Tree

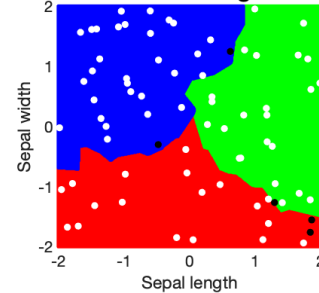


Accuracy = 0.86667

16	6	0	72.7%
21.3%	8.0%	0.0%	27.3%
1	21	1	91.3%
1.3%	28.0%	1.3%	8.7%
2	0	28	93.3%
2.7%	0.0%	37.3%	6.7%
84.2%	77.8%	96.6%	86.7%
15.8%	22.2%	3.4%	13.3%

Target: precision (PPV) and FDR

k-Nearest Neighbor

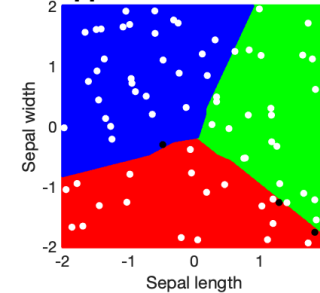


Accuracy = 0.93333

18	3	0	85.7%
24.0%	4.0%	0.0%	14.3%
0	23	0	100%
0.0%	30.7%	0.0%	0.0%
1	1	29	93.5%
1.3%	1.3%	38.7%	6.5%
94.7%	85.2%	100%	93.3%
5.3%	14.8%	0.0%	6.7%

Target: precision (PPV) and FDR

Support Vector Machines

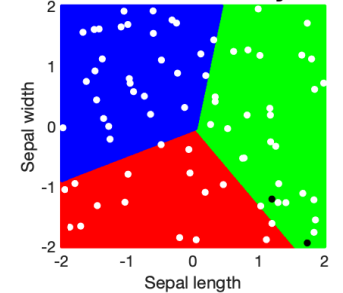


Accuracy = 0.96

18	2	0	90.0%
24.0%	2.7%	0.0%	10.0%
0	25	0	100%
0.0%	33.3%	0.0%	0.0%
1	0	29	96.7%
1.3%	0.0%	38.7%	3.3%
94.7%	92.6%	100%	96.0%
5.3%	7.4%	0.0%	4.0%

Target: precision (PPV) and FDR

Discriminant Analysis



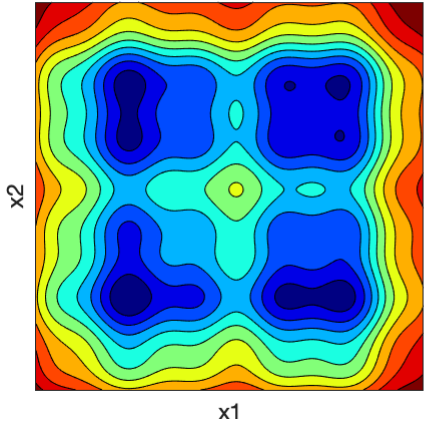
Accuracy = 0.97333

17	0	0	100%
22.7%	0.0%	0.0%	0.0%
2	27	0	93.1%
2.7%	36.0%	0.0%	6.9%
0	0	29	100%
0.0%	0.0%	38.7%	0.0%
89.5%	100%	100%	97.3%
10.5%	0.0%	0.0%	2.7%

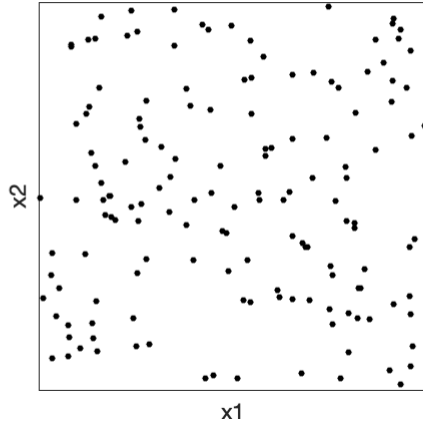
Target: precision (PPV) and FDR



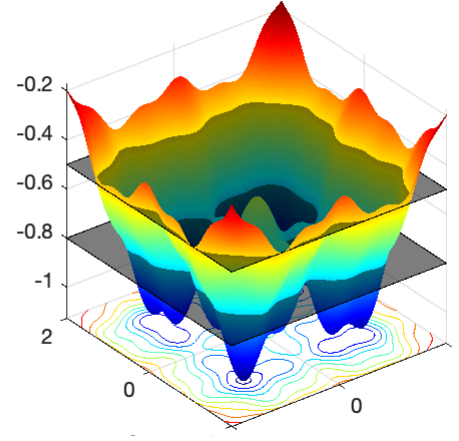
Underlying Behavior



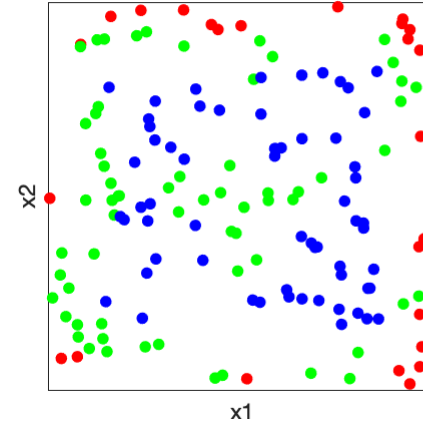
Training Samples (MC)



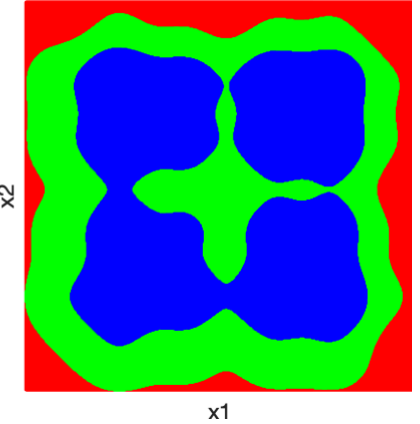
Labeling



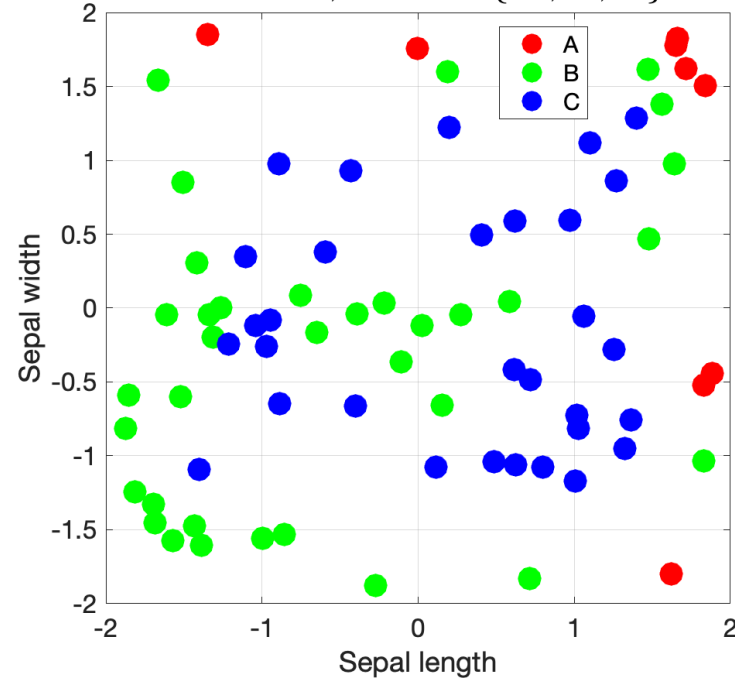
Ground Truth Labels



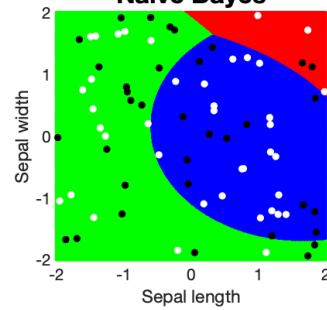
Ideal Classifier



Train: 75, Test: 75{17,28,30}



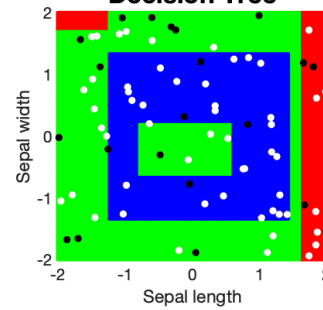
Naive Bayes



Accuracy = 0.53333

Predicted: recall (TPR) and FNR		Target: precision (PPV) and FDR			
	3	2	0		
3	4.0%	2.7%	0.0%	60.0%	40.0%
11	14.7%	21.3%	12.0%	44.4%	55.6%
3	4.0%	13.3%	28.0%	61.8%	38.2%
17.6%	57.1%	70.0%	53.3%	82.4%	46.7%

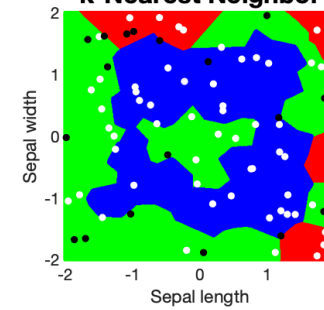
Decision Tree



Accuracy = 0.73333

Predicted: recall (TPR) and FNR		Target: precision (PPV) and FDR			
	7	3	0		
7	9.3%	4.0%	0.0%	70.0%	30.0%
10	13.3%	28.0%	4.0%	61.8%	38.2%
0	0.0%	5.3%	36.0%	87.1%	12.9%
41.2%	75.0%	90.0%	73.3%	58.8%	26.7%

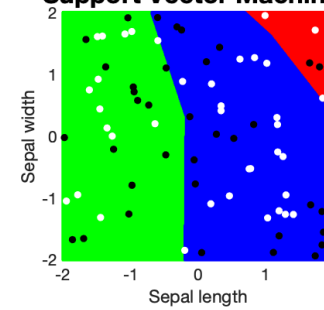
k-Nearest Neighbor



Accuracy = 0.74667

Predicted: recall (TPR) and FNR		Target: precision (PPV) and FDR			
	8	5	0		
8	10.7%	6.7%	0.0%	61.5%	38.5%
9	12.0%	29.3%	5.3%	62.9%	37.1%
0	0.0%	1.3%	34.7%	96.3%	3.7%
47.1%	78.6%	86.7%	74.7%	52.9%	25.3%

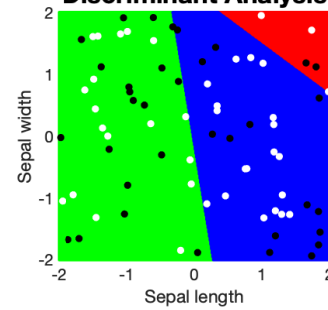
Support Vector Machines



Accuracy = 0.52

Predicted: recall (TPR) and FNR		Target: precision (PPV) and FDR			
	4	2	0		
4	5.3%	2.7%	0.0%	66.7%	33.3%
5	6.7%	20.0%	13.3%	50.0%	50.0%
8	10.7%	14.7%	26.7%	51.3%	48.7%
23.5%	53.6%	66.7%	52.0%	76.5%	48.0%

Discriminant Analysis



Accuracy = 0.53333

Predicted: recall (TPR) and FNR		Target: precision (PPV) and FDR			
	3	2	0		
3	4.0%	2.7%	0.0%	60.0%	40.0%
8	10.7%	24.0%	14.7%	48.6%	51.4%
6	8.0%	10.7%	25.3%	57.6%	42.4%
17.6%	64.3%	63.3%	53.3%	82.4%	46.7%



Ongoing work

- No one solution fits all. Classes of functions vary in terms of smoothness, oscillation, discontinuities ..
- Methods vary in terms of complexity, accuracy, scalability, and computational cost.

Preliminary Results

- Deep ensembles seem to offer multiple advantages for regression problems, including i) smoothed out performance by adversarial training, ii) conservative estimates, and iii) low training cost (with $M=5$ models)
- Comparisons to hybrid Bayesian/non-Bayesian methods

Thank you!