

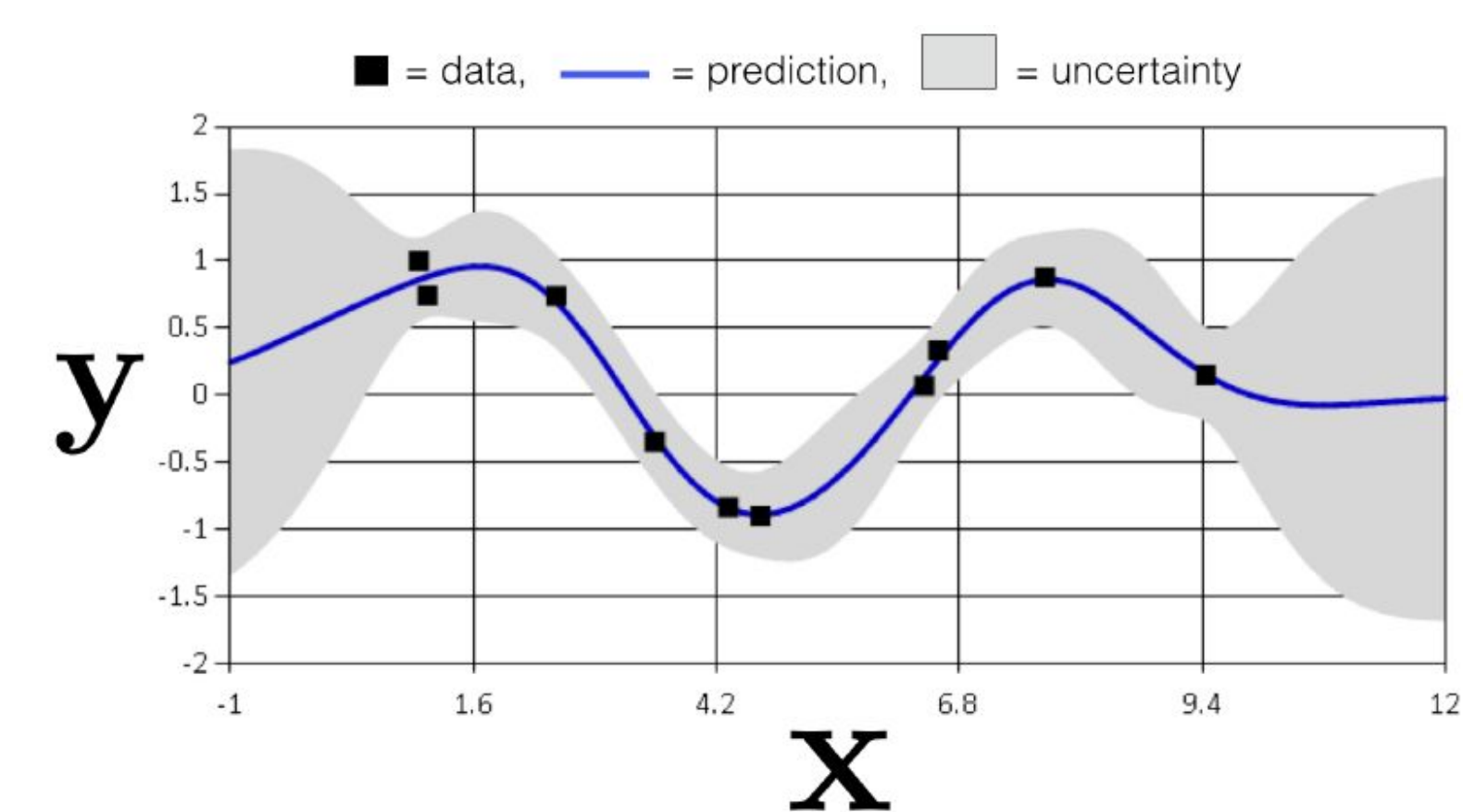
HYBRID MODELS WITH DEEP AND INVERTIBLE FEATURES

Eric Nalisnick*, Akihiro Matsukawa*, Yee Whye Teh, Dilan Gorur, Balaji Lakshminarayanan

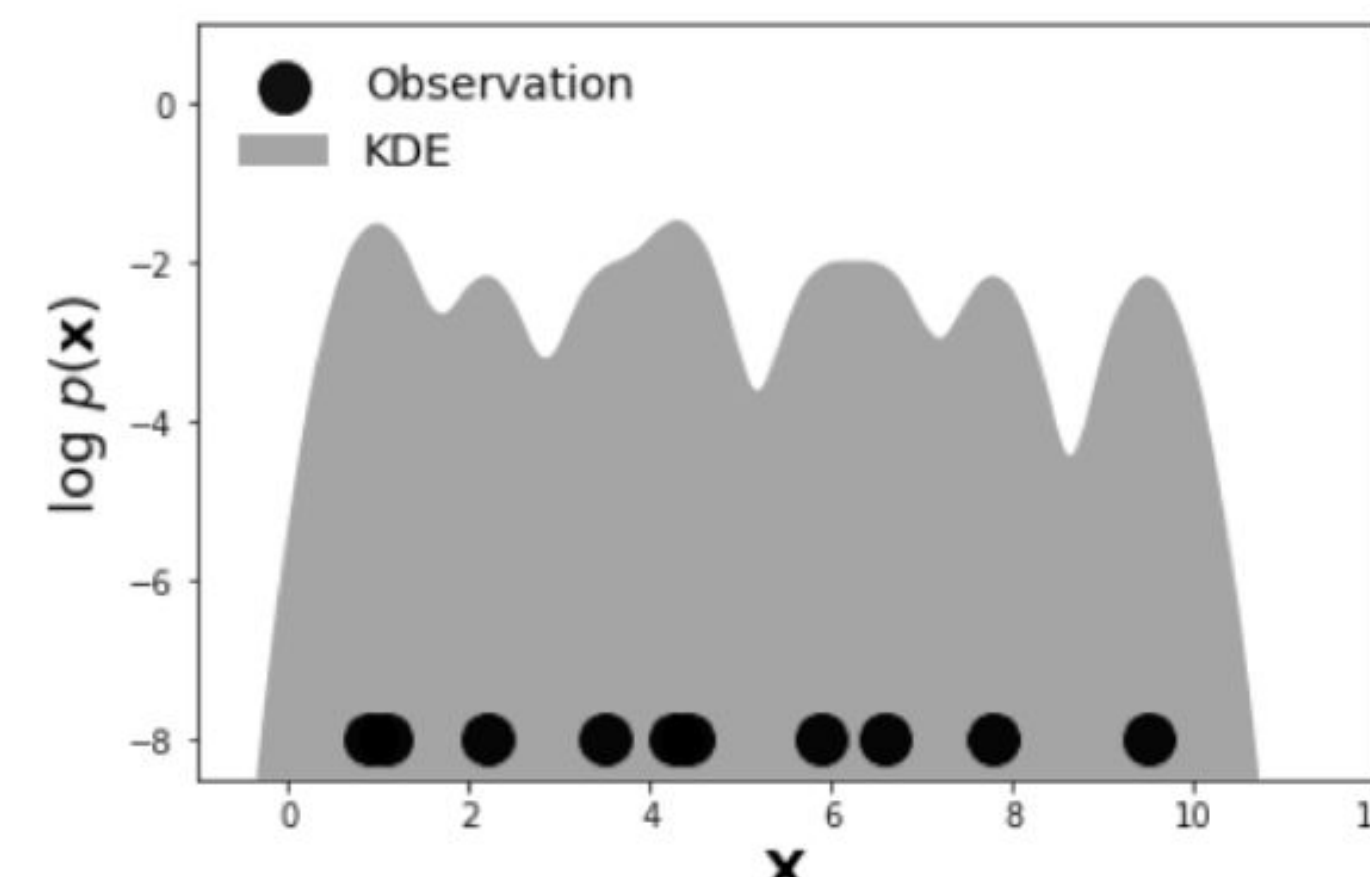
1. INTRODUCTION

- Neural networks usually model the conditional distribution $p(y|x)$, where y denotes a label and x features.
- Generative models, on the other hand, represent the distribution over features $p(x)$.
- Can we efficiently combine the two in a hybrid model of the joint distribution $p(y, x)$?

Conditional Model



Generative Model



2. BACKGROUND

Invertible Generative Models (Normalizing Flows)

Invertible generative models (a.k.a. normalizing flows) are a broad class of models defined via the change-of-variables formula. An initial density $p(x)$ 'flows' through a series of transformations $f(x)$ and morphs into some (usually simpler) prior distribution $p(z)$.

$$\log p_x(x) = \log p_z(f(x; \phi)) + \log \left| \frac{\partial f_\phi}{\partial x} \right|$$

Generalized Linear Models (GLMs)

Generalized linear models (GLMs) model the expected response (or label) y as a transformation of the linear model $\beta^T z$ where β are parameters and z are features (covariates).

$$\mathbb{E}[y_n | z_n] = g^{-1}(\beta^T z_n)$$

- Regression:** $\mathbb{E}[y|z] = \text{identity}(\beta^T z)$
- Binary Classification:** $\mathbb{E}[y|z] = \text{logistic}(\beta^T z)$

3. COMBINING DEEP GENERATIVE MODELS AND LINEAR MODELS

We define a model of the joint distribution $p(y, x)$ by instantiating a GLM on the output of a normalizing flow:

$$p(y_n, x_n; \theta) = p(y_n | x_n; \beta, \phi) p(x_n; \phi)$$

$$= p(y_n | f(x_n; \phi); \beta) p_z(f(x_n; \phi)) \left| \frac{\partial f_\phi}{\partial x_n} \right|$$

In practice, we add a weight to the flow terms to tradeoff between predictive and generative behavior:

$$\mathcal{J}_\lambda(\theta) = \sum_{n=1}^N \left(\log p(y_n | x_n; \beta, \phi) + \lambda \log p(x_n; \phi) \right)$$

Examples

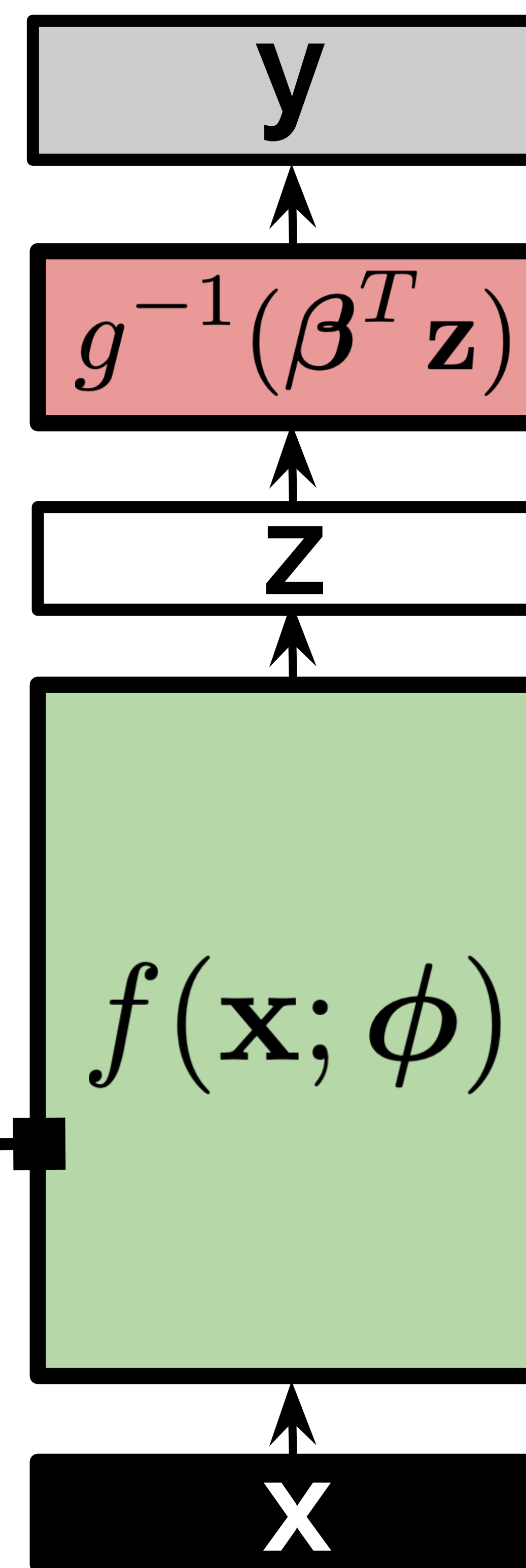
Planar: $= |1 + \mathbf{u}^T f'(\mathbf{w}^T \mathbf{x} + b) \mathbf{w}|$ where \mathbf{w}, \mathbf{u} are parameters.
RNVP: $= \sum_l \sum_d s_{l,d}(\mathbf{x}; \phi)$ where $s(l)$ are scaling operations.
Glow: $= \sum_l \sum_d s_{l,d}(\mathbf{x}; \phi) + h_l w_l \log |\det \mathbf{W}_l|$, $\mathbf{w}_{1 \times 1}$ params.

Bayesian treatment: we can place a prior on the parameters of the GLM in order to quantify model and data uncertainty.

$$f(x; \phi) \sim p(z), \quad \beta \sim p(\beta), \quad y_n \sim p(y_n | f(x_n; \phi), \beta)$$

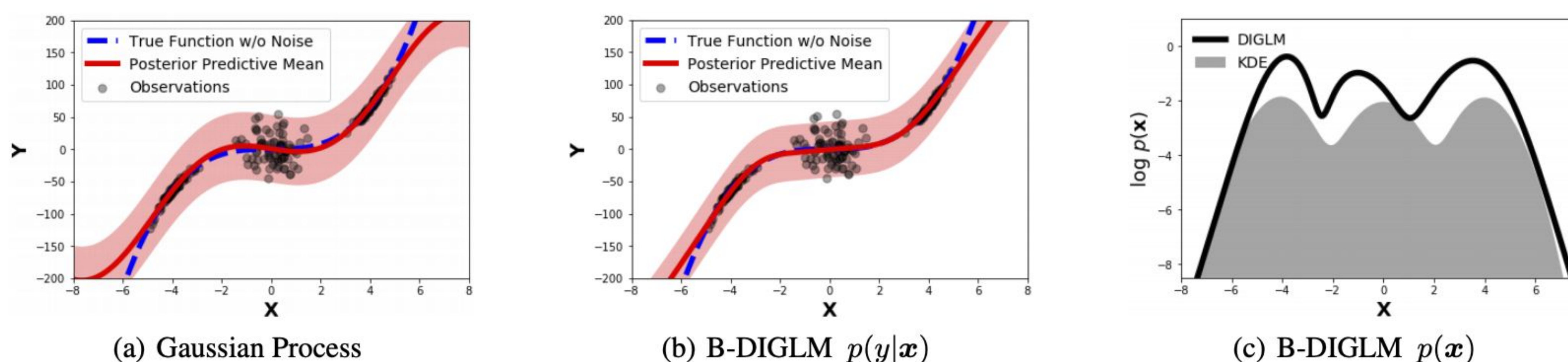
For a Gaussian prior on the GLM, the predictive model can be trained via the closed-form marginal likelihood:

$$\log p(y_n | f(x_n; \phi)) = \log N(y; \mathbf{0}, \sigma_0^2 \mathbb{I} + \lambda^{-1} \mathbf{Z}_\phi \mathbf{Z}_\phi^T)$$



Deep Invertible Generalized Linear Model (DIGLM)

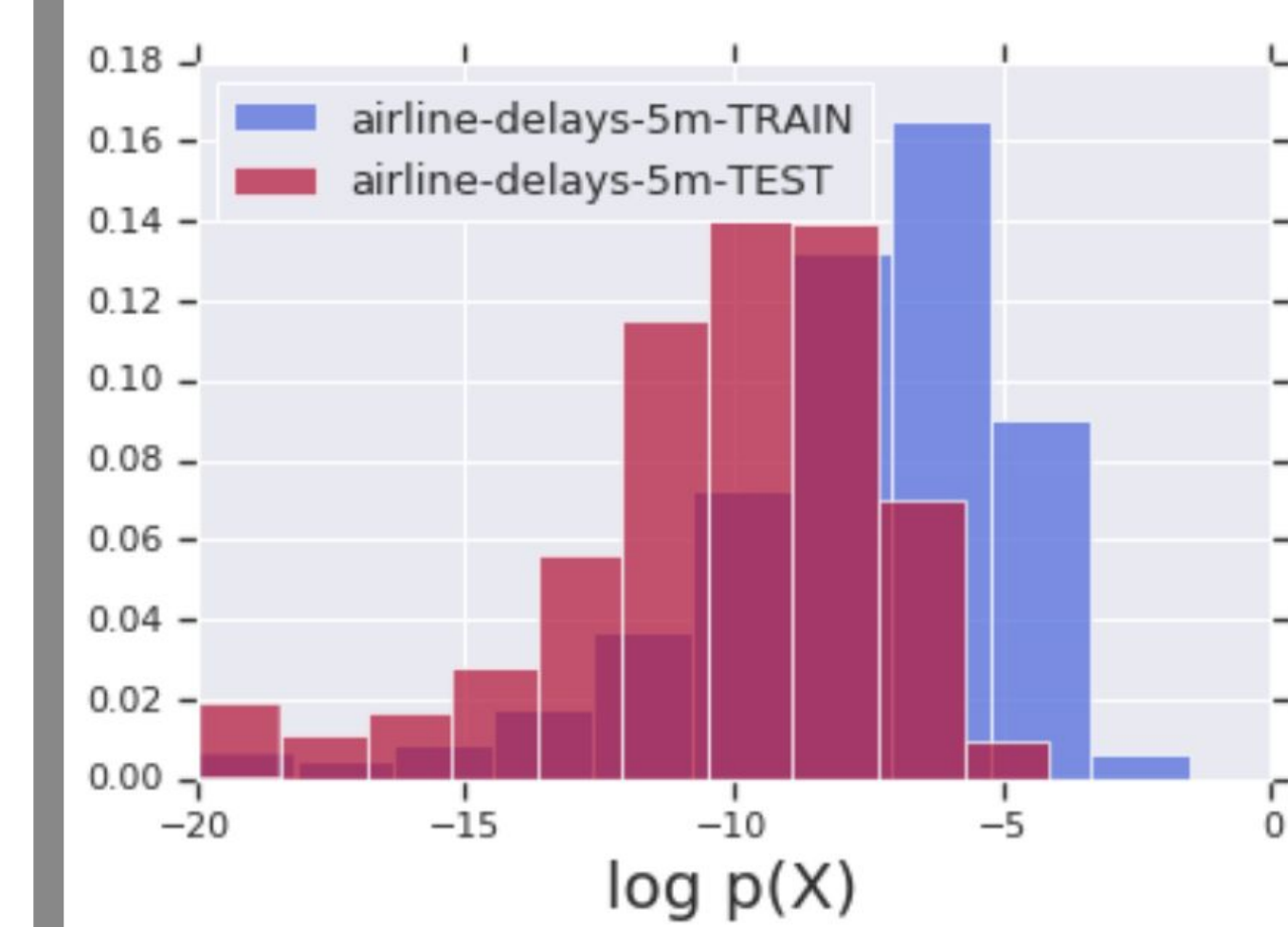
4. SIMULATION



1D regression task with heteroscedastic noise. **Subfigure (a)** shows a Gaussian process and **Subfigure (b)** shows our Bayesian DIGLM. **Subfigure (c)** shows $p(x)$ learned by the same DIGLM (black line) and compares it to a KDE (gray shading).

5. EXPERIMENTS

Regression on Flight Delay Data Set (N=5 million, D=8)

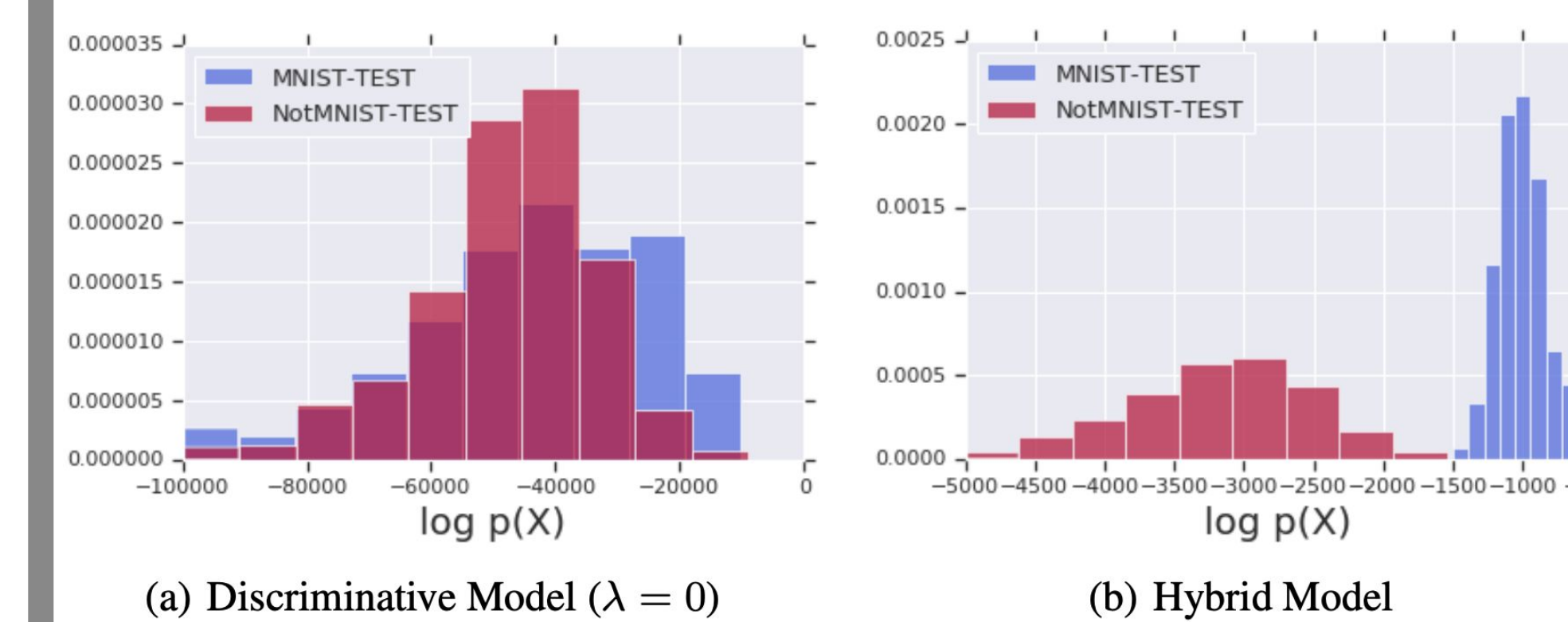


MODEL	RMSE ↓	NLL ↓
MONDRIAN FORESTS (SOTA)	38.38	6.91
DIGLM	40.46	5.07

- This data set exhibits covariate shift between the train and test splits.
- The DIGLM's $p(x)$ component is able to detect this shift (see left).

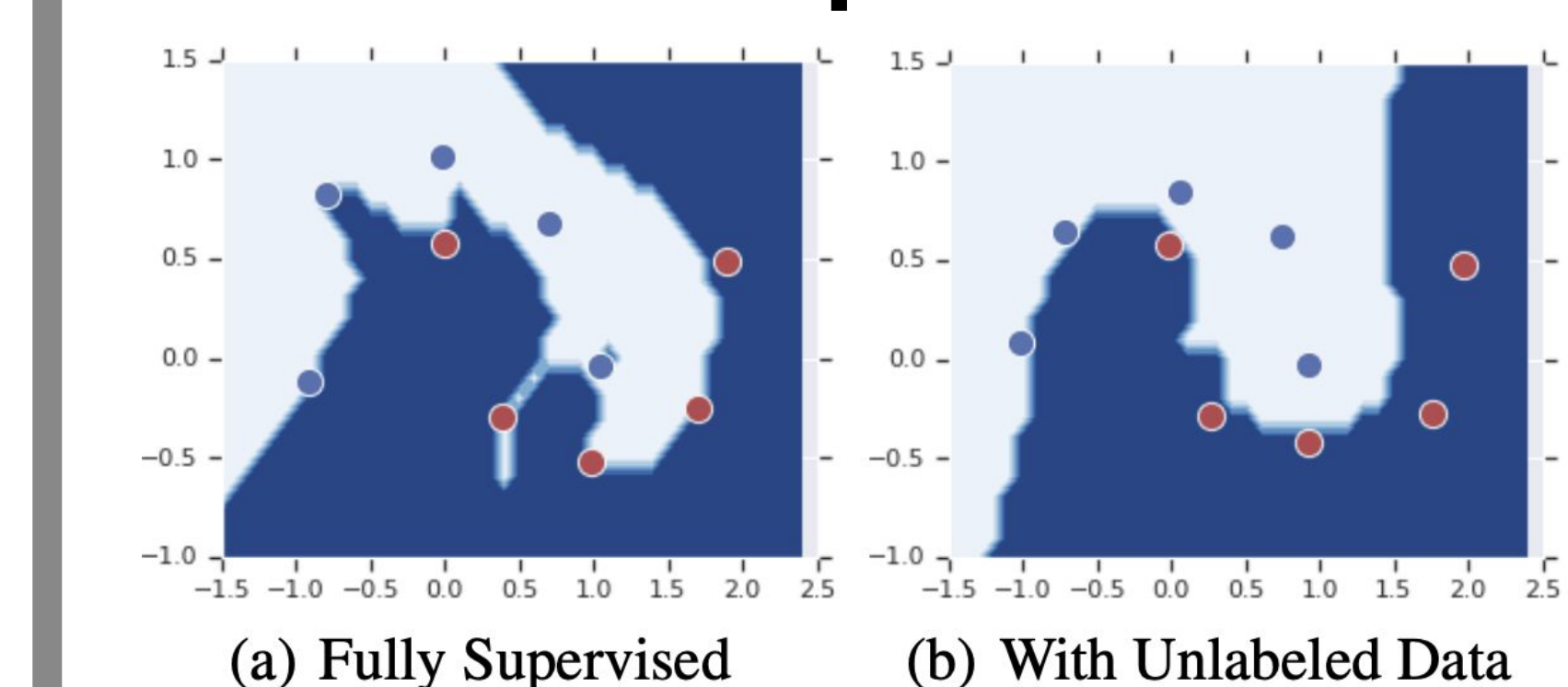
Classification on MNIST and SVHN

Model	BPD ↓	MNIST			NotMNIST			Model	BPD ↓	SVHN			CIFAR-10		
		error ↓	NLL ↓	Entropy ↑	BPD ↓	NLL ↓	Entropy ↑			BPD ↓	NLL ↓	Entropy ↑			
Discriminative ($\lambda = 0$)	81.80*	0.67%	0.082	87.74*	29.27	0.130	Discriminative ($\lambda = 0$)	15.40*	4.26%	0.225	15.20*	4.60	0.998		
Hybrid ($\lambda = 0.01/D$)	1.83	0.73%	0.035	5.84	2.36	2.300	Hybrid ($\lambda = 0.1/D$)	3.35	4.86%	0.260	7.06	5.06	1.153		
Hybrid ($\lambda = 1.0/D$)	1.26	2.22%	0.081	6.13	2.30	2.300	Hybrid ($\lambda = 1.0/D$)	2.40	5.23%	0.253	6.16	4.23	1.677		
Hybrid ($\lambda = 10.0/D$)	1.25	4.01%	0.145	6.17	2.30	2.300	Hybrid ($\lambda = 10.0/D$)	2.23	7.27%	0.268	7.03	2.69	2.143		



- λ controls the trade-off between $p(y|x)$ and $p(x)$.
- Hybrid model is better able to detect the OOD inputs via $p(x)$.

Semi-Supervised Learning: MNIST and Half Moons



Half-moons simulation: the DIGLM leverages unlabeled data to learn a smooth decision boundary (N=10 labeled points).

Model	MNIST-error ↓	MNIST-NLL ↓	SSL (VAT) with only 1000 labels (2% of labeled data) achieves <1% error on MNIST
1000 labels only	6.61%	0.276	
1000 labels + unlabeled	0.99%	0.069	
All labeled	0.73%	0.035	

6. SUMMARY

We defined a neural hybrid model that can efficiently compute both predictive $p(y|x)$ and generative $p(x)$ distributions, in a single feed-forward pass, making it a useful building block for downstream applications of probabilistic deep learning.

Paper: <https://arxiv.org/abs/1902.02767>