

Hybrid Models with Deep and Invertible Features

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

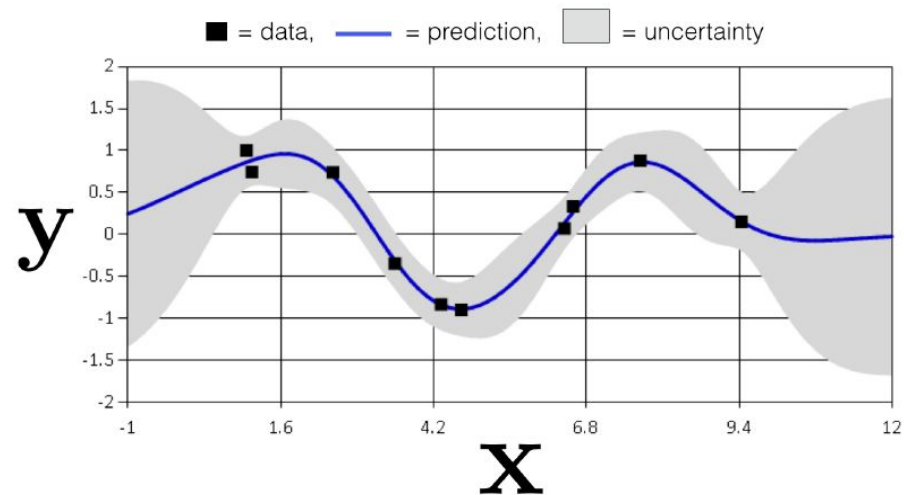


DeepMind

*equal contribution

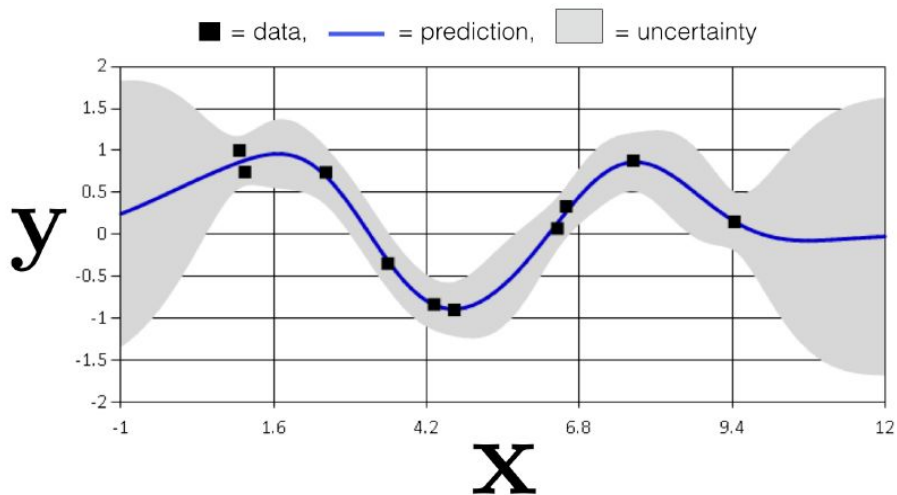
Predictive Models

$$p(\mathbf{y} | \mathbf{x}; \boldsymbol{\theta})$$



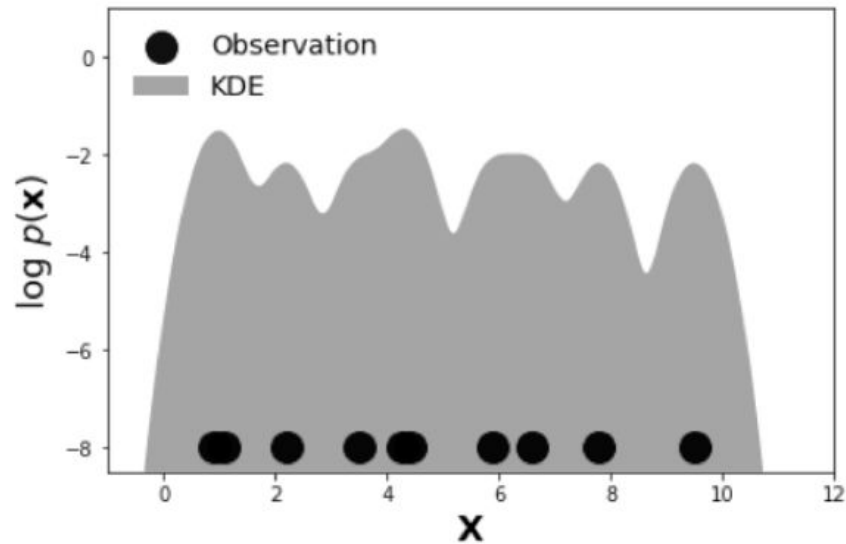
Predictive Models

$$p(\mathbf{y}|\mathbf{x}; \boldsymbol{\theta})$$



Generative Models

$$p(\mathbf{x}; \phi)$$



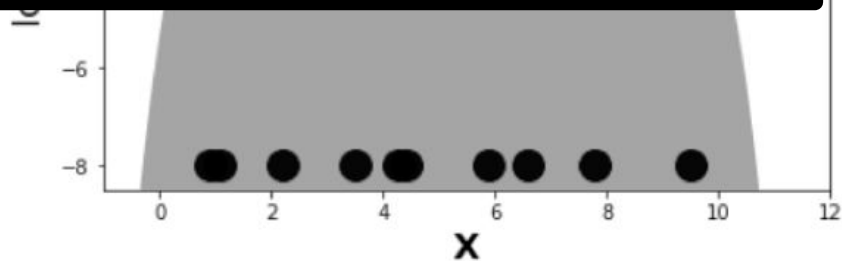
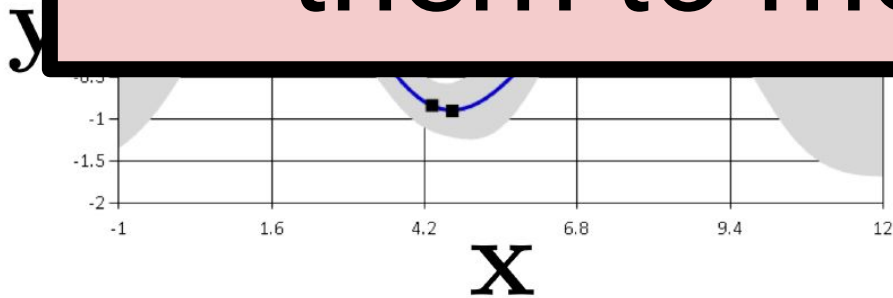
Predictive Models

$$p(\mathbf{y} | \mathbf{x}; \boldsymbol{\theta})$$

Generative Models

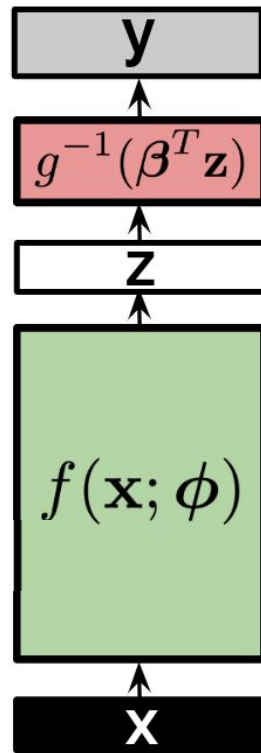
$$p(\mathbf{x}; \boldsymbol{\phi})$$

Can we efficiently combine them to model $p(y, x)$?



Neural Hybrid Model

We define a computationally efficient **hybrid model** by combining *normalizing flows* with *generalized linear models* (GLMs).

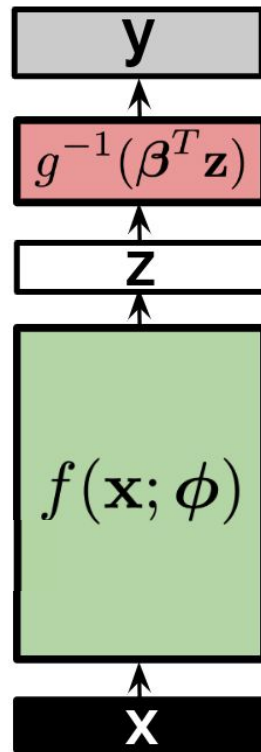


Deep Invertible
Generalized
Linear Model
(DIGLM)

Neural Hybrid Model

We define a computationally efficient **hybrid model** by combining *normalizing flows* with *generalized linear models* (GLMs).

$$p(y_n, \mathbf{x}_n; \theta) = \underbrace{p(y_n | \mathbf{x}_n; \beta, \phi)}_{\text{Predictive Component}} \underbrace{p(\mathbf{x}_n; \phi)}_{\text{Generative Component}}$$



Deep Invertible
Generalized
Linear Model
(DIGLM)

Neural Hybrid Model

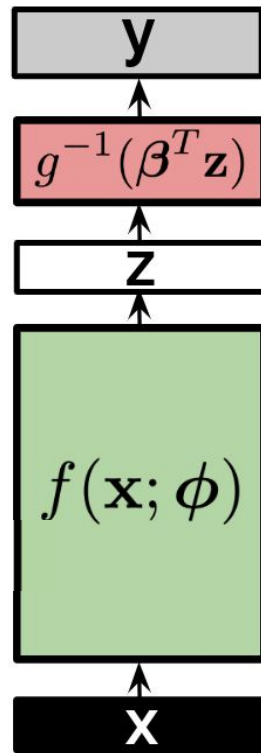
We define a computationally efficient **hybrid model** by combining *normalizing flows* with *generalized linear models* (GLMs).

$$p(y_n, \mathbf{x}_n; \boldsymbol{\theta}) = p(y_n | \mathbf{x}_n; \boldsymbol{\beta}, \phi) p(\mathbf{x}_n; \phi)$$

$$= \boxed{p(y_n | f(\mathbf{x}_n; \phi); \boldsymbol{\beta})} \quad \boxed{p_z(f(\mathbf{x}_n; \phi)) \quad \left| \quad \frac{\partial f_\phi}{\partial \mathbf{x}_n} \right|}$$

Linear Model

Normalizing Flow

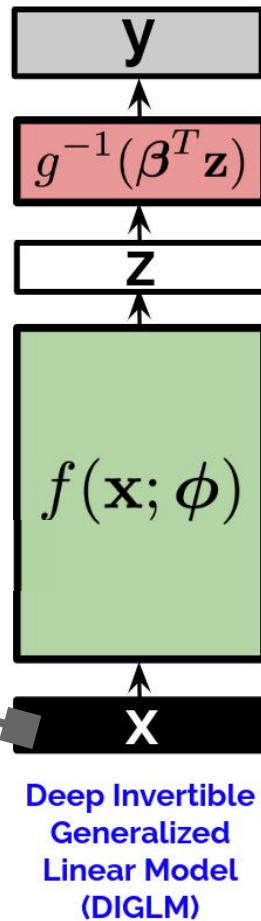


Deep Invertible
Generalized
Linear Model
(DIGLM)

Neural Hybrid Model

We define a computationally efficient **hybrid model** by combining *normalizing flows* with *generalized linear models* (GLMs).

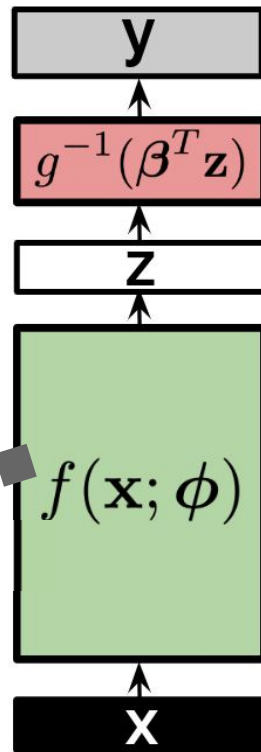
Input features.



Neural Hybrid Model

We define a computationally efficient **hybrid model** by combining *normalizing flows* with *generalized linear models* (GLMs).

Normalizing flow acts as a deep neural feature extractor.

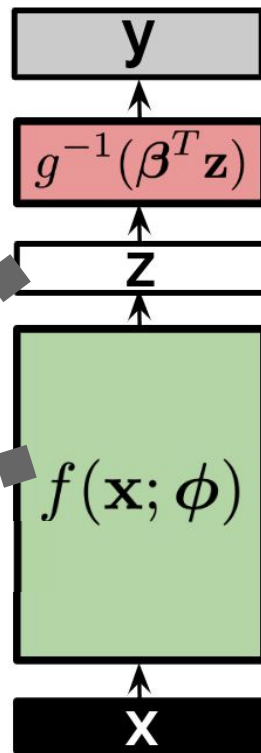


Deep Invertible
Generalized
Linear Model
(DIGLM)

Neural Hybrid Model

We define a computationally efficient **hybrid model** by combining *normalizing flows* with *generalized linear models* (GLMs).

Flow's output and params. are used to compute $p(\mathbf{x})$ via change-of-variables.

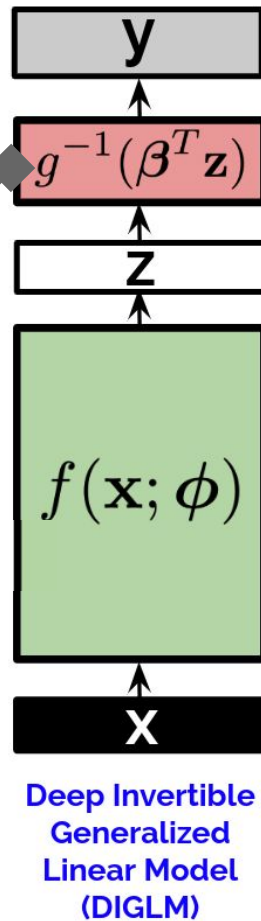


Deep Invertible
Generalized
Linear Model
(DIGLM)

Neural Hybrid Model

We define a computationally efficient **hybrid model** by combining *normalizing flows* with *generalized linear models* (GLMs).

Flow's output is used as the feature vector in a (generalized) linear model, which computes $p(\mathbf{y}|\mathbf{x})$.



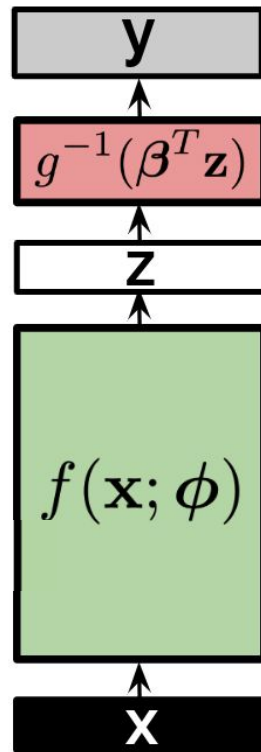
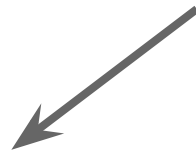
Neural Hybrid Model

We define a computationally efficient **hybrid model** by combining *normalizing flows* with *generalized linear models* (GLMs).

Optimization objective:

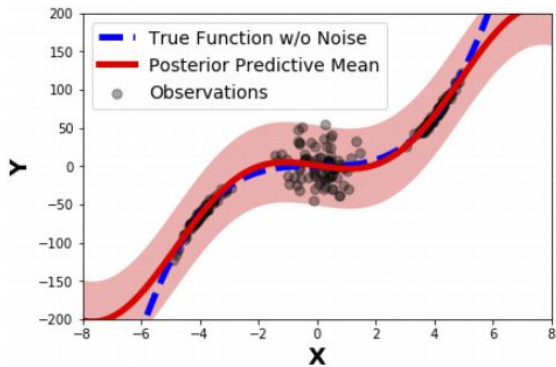
$$\mathcal{J}_\lambda(\theta) = \sum_{n=1}^N \left(\log p(y_n | \mathbf{x}_n; \beta, \phi) + \lambda \log p(\mathbf{x}_n; \phi) \right)$$

Weight to trade-off predictive and generative performance.



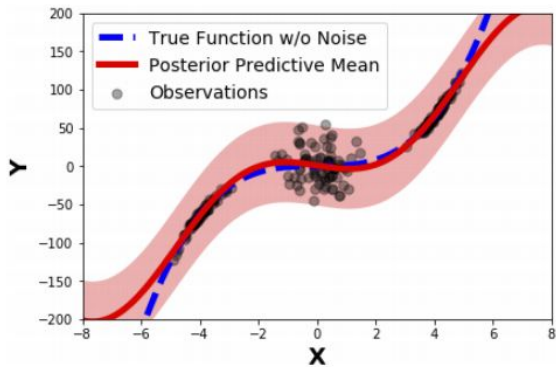
Deep Invertible
Generalized
Linear Model
(DIGLM)

Simulation: Heteroscedastic Regression

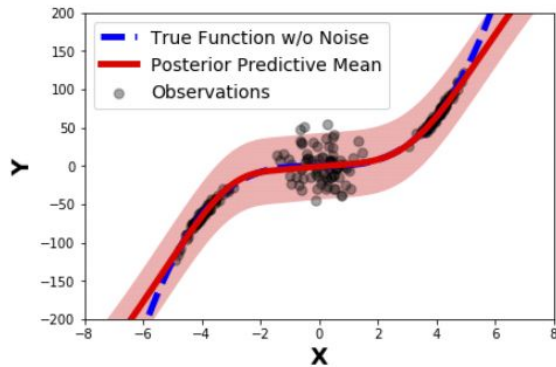


*Gaussian
process fitted to
simulated data.*

Simulation: Heteroscedastic Regression

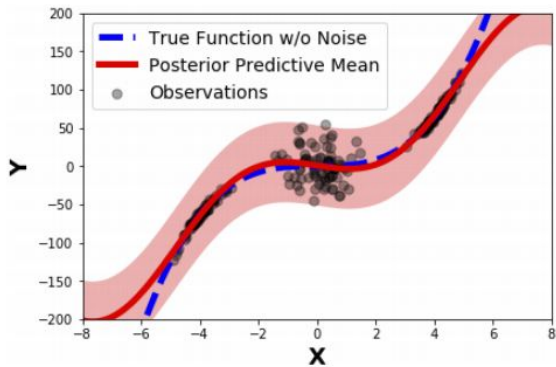


*Gaussian
process fitted to
simulated data.*

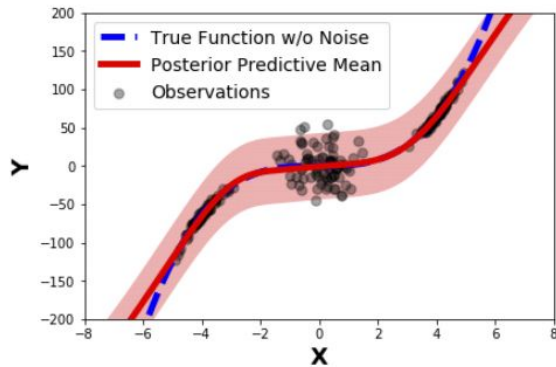


*Our model's
predictive
component.*

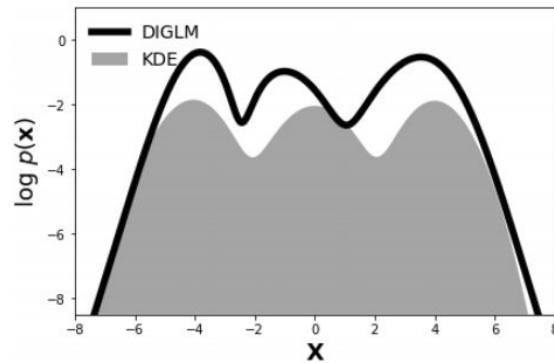
Simulation: Heteroscedastic Regression



Gaussian process fitted to simulated data.



Our model's predictive component.



Our model's generative component.

For more details, please visit our poster.

HYBRID MODELS WITH DEEP AND INVERTIBLE FEATURES

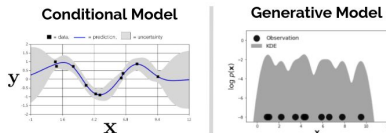


Eric Nalisnick*, Akihiro Matsukawa*, Yee Whye Teh, Dilan Gorur, Balaji Lakshminarayanan

*equal contribution

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 $p(y, x)$?



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 $f(x; \phi)$ 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{u} are parameters.

RNVP: $= \sum_i \sum_d s_{i,d}(\mathbf{x}; \phi)$ where s_i are scaling operations.

Glow: $= \sum_i \sum_d s_{i,d}(\mathbf{x}; \phi) + h_i w_i \log |\det \mathbf{W}_i|$, \mathbf{W}_i x1x1 params

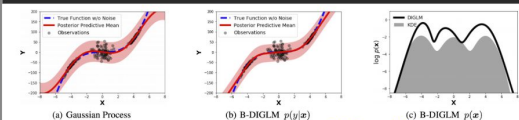
Bayesian treatment: we can place a prior on the parameters of the GLM in order to quantify model and data uncertainty.

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

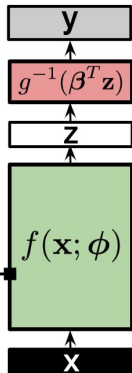
For a Gaussian prior on the GLM, the predictive model can be trained via the closed-form marginal likelihood:

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

4. SIMULATION



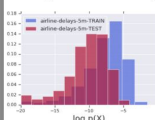
1D regression task with heteroscedastic noise. **Subfigure (a)** shows a Gaussian process and **Subfigure (b)** shows our Bayesian DIGLM. **Subfigure (c)** shows $p(x)$ learned by the same DIGLM (black line) and compares it to a KDE (gray shading).



Deep Invertible Generalized Linear Model (DIGLM)

5. EXPERIMENTS

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

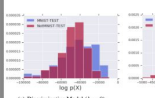


MODEL	RMSE ↓	NLL ↓
MONDRIAN FORESTS (SOTA)	38.38	6.91
DIGLM	40.46	5.07

- 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).

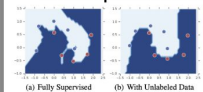
Classification on MNIST and SVHN

Model	MNIST			SVHN			CIFAR-10						
	RFD ↓	error ↓	NLL ↓	RFD ↓	error ↓	NLL ↓	RFD ↓	error ↓	NLL ↓				
Discriminative ($\lambda = 0$)	81.80	6.47%	0.082	87.54	29.27	0.110	15.60	4.26%	0.228	15.20	4.60	0.996	
Hybrid ($\lambda = 0.01$)	1.83	0.75%	0.009	5.94	2.36	2.300	Hybrid ($\lambda = 0.1$)	3.35	4.65%	0.261	7.06	5.06	1.153
Hybrid ($\lambda = 1.0$)	1.26	2.25%	0.081	6.13	2.30	2.300	Hybrid ($\lambda = 1.0$)	2.40	5.21%	0.251	6.16	4.33	1.047
Hybrid ($\lambda = 10.0$)	1.28	6.01%	0.145	6.47	2.30	2.300	Hybrid ($\lambda = 10.0$)	2.23	7.21%	0.268	2.68	2.69	0.848



- 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>