
Supplementary Materials: Hybrid Models with Deep and Invertible Features

Eric Nalisnick^{*1} Akihiro Matsukawa^{*1} Yee Whye Teh¹ Dilan Gorur¹ Balaji Lakshminarayanan¹

A. Background on Affine Coupling Layers

One ACL performs the following operations (Dinh et al., 2017):

1. *Splitting*: \mathbf{x} split it at dimension d into two separate vectors $\mathbf{x}_{:d}$ and \mathbf{x}_d . (using Python list syntax).
2. *Identity and Affine Transformations*: Given the split $\{\mathbf{x}_{:d}, \mathbf{x}_d\}$, each undergoes a separate operation:

$$\text{identity} : \mathbf{h}_{:d} = \mathbf{x}_{:d} \quad (1)$$

$$\text{affine} : \mathbf{h}_d = t(\mathbf{x}_{:d}; \phi_t) + \mathbf{x}_d \odot \exp\{s(\mathbf{x}_{:d}; \phi_s)\}$$

where $t(\cdot)$ and $s(\cdot)$ are translation and scaling operations with no restrictions on their functional form. We can compute them with neural networks that take as input $\mathbf{x}_{:d}$, the other half of the original vector, and since $\mathbf{x}_{:d}$ has been copied forward by the first operation, no information is lost that would jeopardize invertibility.

3. *Permutation*: Lastly, the new representation $\mathbf{h} = \{\mathbf{h}_{:d}, \mathbf{h}_d\}$ is ready to be either treated as output or fed into another ACL. If the latter, then the elements should be modified so that $\mathbf{h}_{:d}$ is not again copied but rather subject to the affine transformation. Dinh et al. (2017) simply exchange the components (i.e. $\{\mathbf{h}_{:d}, \mathbf{h}_d\}$) whereas Kingma & Dhariwal (2018) apply a 1×1 convolution, which can be thought of as a continuous generalization of a permutation.

Several ACLs are composed to create the final form of $f(\mathbf{x}; \phi)$, which is called a *normalizing flow* (Rezende & Mohamed, 2015). Crucially, the Jacobian of these operations is efficient to compute, simplifying to the sum of all the scale transformations:

$$\begin{aligned} \log \left| \frac{\partial \mathbf{f}_\phi}{\partial \mathbf{x}} \right| &= \log \exp \left\{ \sum_{l=1}^L s_l(\mathbf{x}_{:d}; \phi_{s,l}) \right\} \\ &= \sum_{l=1}^L s_l(\mathbf{x}_{:d}; \phi_{s,l}) \end{aligned}$$

where l is an index over ACLs. The Jacobian of a $1 \times$

^{*}Equal contribution ¹DeepMind. Correspondence to: Balaji Lakshminarayanan <balajiln@google.com>.

1 convolution does not have as simple of an expression, but Kingma & Dhariwal (2018) describe ways to reduce computation.

B. Additional Semi-Supervised Results

B.1. Hyperparameters for Semi-supervised learning

The hyperparameters are described in Table S1.

Hyper-parameter	Grid values
Dropout rate	0, 0.2, 0.5
ϵ_{VAT} for Virtual Adversarial Training	1, 5
λ_{EM} for Entropy Minimization loss	0, 0.3, 1, 3
λ_{VAT} for Virtual Adversarial Training loss	0, 0.3, 1, 3

Table S1. Hyperparameters for the semi-supervised learning experiments.

B.2. Semi-supervised learning on SVHN

The results are shown in Table S2. Similar to the semi-supervised results on MNIST, we observe that our model can effectively leverage unlabeled data.

Model	SVHN-Error ↓	SVHN-NLL ↓
1000 labels only	19.26%	0.78
1000 labels + unlabeled	5.90%	0.38
All labeled	4.86%	0.26
All labeled + unlabeled	2.80%	0.17

Table S2. Results of hybrid model for semi-supervised learning on SVHN. Arrows indicate which direction is better.

C. Extensions

C.1. Mixed Effects Model

To model heteroscedastic noise, we can also add “random effects” to the model at the latent level. The model is then:

$$\begin{aligned} \mathbb{E}[y_n | \mathbf{x}_n] &= g^{-1} \left(\beta^T f(\mathbf{x}_n; \phi) + \mathbf{u}^T \mathbf{a}_n \right), \\ \mathbf{u} &\sim p(\mathbf{u}), \quad \mathbf{a}_n \sim p(\mathbf{a}) \end{aligned}$$

where \mathbf{a}_n is a vector of random effects associated with the fixed effects \mathbf{x}_n , and \mathbf{u} are the corresponding parameters in the GLM. Note that Depeweg et al. (2018) also use random effects model to handle heteroscedastic noise, but add them at the input level.

References

- Depeweg, S., Hernández-Lobato, J. M., Doshi-Velez, F., and Udluft, S. Decomposition of Uncertainty in Bayesian Deep Learning for Efficient and Risk-sensitive Learning. In *ICML*, 2018.
- Dinh, L., Sohl-Dickstein, J., and Bengio, S. Density Estimation Using Real NVP. In *International Conference on Learning Representations (ICLR)*, 2017.
- Kingma, D. P. and Dhariwal, P. Glow: Generative Flow with Invertible 1x1 Convolutions. In *Advances in Neural Information Processing Systems (NeurIPS)*, 2018.
- Rezende, D. and Mohamed, S. Variational Inference with Normalizing Flows. In *International Conference on Machine Learning (ICML)*, 2015.