# Deep Probabilistic Programming with Edward

Dustin Tran[†], Matt Hoffman[*‡], Kevin Murphy[‡], Eugene Brevdo[‡], Rif Saurous[‡], David Blei[†]

[†]Columbia University, [*]Adobe Research, [‡]Google

edwardlib.org/iclr2017

## Summary

- Deep neural networks are popular in large part due to their compositional nature. How do we do this for probabilistic modeling?
- We describe Edward, a Turing-complete probabilistic programming language.
- Edward builds two representations—random variables and inference.
- For example, we show how to design rich variational models and generative adversarial networks.

## Compositional Representations for Probabilistic Models

- We define random variables as the key compositional representation.
- They are class objects e.g. with log-density and sample methods.
- Each random variable $\mathbf{x}$ is associated to a tensor $\mathbf{x}^*$ in the computational graph, which represents a single sample $\mathbf{x}^* \sim p(\mathbf{x})$.
- Mutable states represent enable conditioning sets to vary, $p(\mathbf{y}|\mathbf{x})$ and optimization of parameters, $p(\mathbf{x}; \theta)$.

## Compositional Representations for Inference

- Given data $\mathbf{x}_{\text{train}}$, inference aims to calculate the posterior $p(\mathbf{z}, \beta \mid \mathbf{x}_{\text{train}}; \theta)$, where $\theta$ are any model parameters to estimate.
- In variational inference, the idea is to posit an approximating family $q \in \mathscr{Q}$ and to find the closest member $q^*$. We write it with mutable states representing its parameters, where $q(\beta; \mu, \sigma) = \text{Normal}(\beta; \mu, \sigma)$, $q(\mathbf{z}; \pi) = \text{Categorical}(\mathbf{z}; \pi)$.

```
1  qbeta = Normal(mu=tf.Variable(tf.zeros([K, D])),
2          sigma=tf.exp(tf.Variable(tf.zeros[K, D])))
3  qz = Categorical(logits=tf.Variable(tf.zeros[N, K]))
4
5  inference = ed.VariationalInference({beta: qbeta, z: qz}, data={x: x_train}
```

- Specific variational algorithms inherit from VariationalInference to define their own methods, e.g., a loss function and gradient.
- Monte Carlo approximates the posterior using samples. We represent it where the approximating family is an empirical distribution, $q(\beta; \{\beta^{(t)}\}) = \frac{1}{T}\sum_{t=1}^{T} \delta(\beta, \beta^{(t)})$, $q(\mathbf{z}; \{\mathbf{z}^{(t)}\}) = \frac{1}{T}\sum_{t=1}^{T} \delta(\mathbf{z}, \mathbf{z}^{(t)})$.
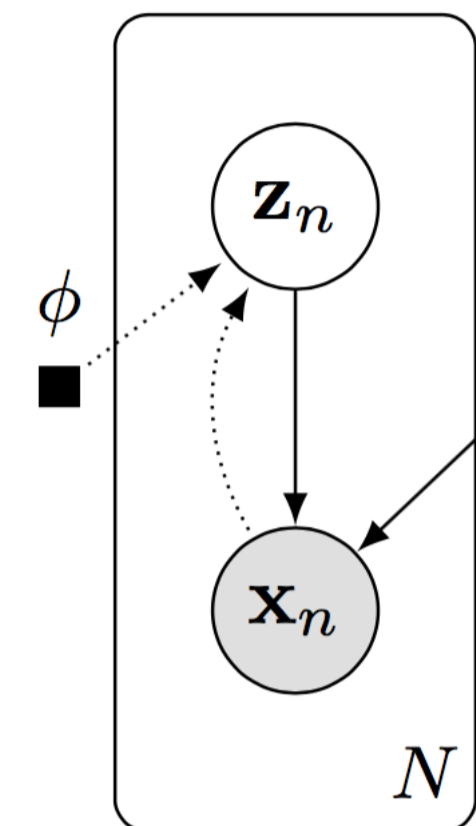
```
1  T = 10000  # number of samples
2  qbeta = Empirical(params=tf.Variable(tf.zeros([T, K, D]))
3  qz = Empirical(params=tf.Variable(tf.zeros([T, N]))
4
5  inference = ed.MonteCarlo({beta: qbeta, z: qz}, data={x: x_train})
```

- Monte Carlo algorithms proceed by updating one sample $\beta^{(t)}, \mathbf{z}^{(t)}$ at a time in the empirical approximation. Specific MC samplers determine the update rules.

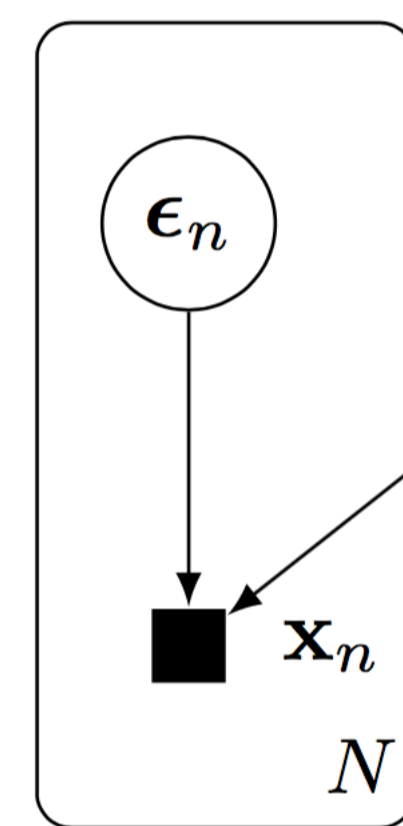## Example: Variational Auto-Encoder



```
1   # Probabilistic model
2   z = Normal(mu=tf.zeros([N, d]), sigma=tf.ones([N, d]))
3   h = Dense(256, activation='relu')(z)
4   x = Bernoulli(logits=Dense(28 * 28, activation=None)(h))
5
6   # Variational model
7   qx = tf.placeholder(tf.float32, [N, 28 * 28])
8   qh = Dense(256, activation='relu')(qx)
9   qz = Normal(mu=Dense(d, activation=None)(qh),
10          sigma=Dense(d, activation='softplus')(qh))
```

## Example: Generative Adversarial Networks
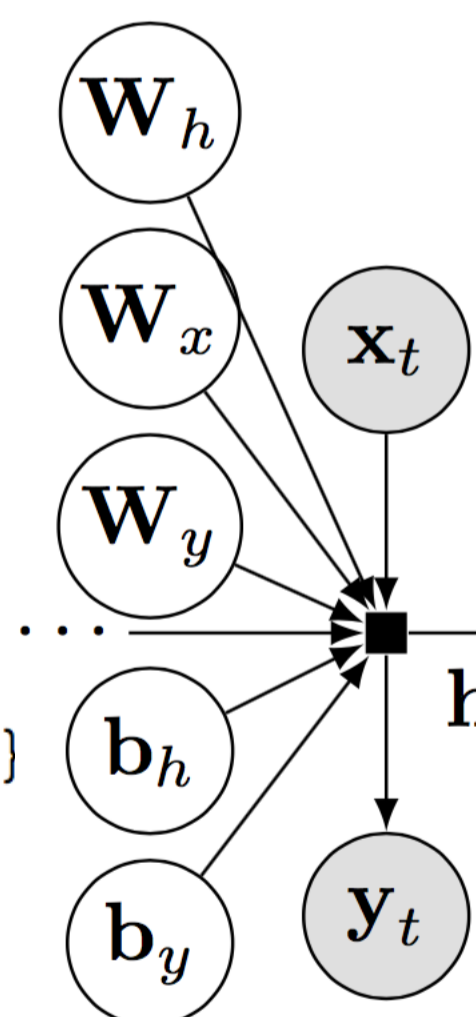


```
1    def generative_network(eps):
2      h = Dense(256, activation='relu')(eps)
3      return Dense(28 * 28, activation=None)(h)
4
5    def discriminative_network(x):
6      h = Dense(28 * 28, activation='relu')(x)
7      return Dense(h, activation=None)(1)
8
9    # Probabilistic model
10   eps = Normal(mu=tf.zeros([M, d]), sigma=tf.ones([M, d]))
11   x = generative_network(eps)
12
13   inference = ed.GANInference(data={x: x_train},
14         discriminator=discriminative_network)
```

## Example: Bayesian RNN with Variable Length



```
1    def rnn_cell(hprev, xt):
2      return tf.tanh(tf.dot(hprev, Wh) + tf.dot(xt, Wx) + bh)
3
4    Wh = Normal(mu=tf.zeros([H, H]), sigma=tf.ones([H, H]))
5    Wx = Normal(mu=tf.zeros([D, H]), sigma=tf.ones([D, H]))
6    Wy = Normal(mu=tf.zeros([H, 1]), sigma=tf.ones([H, 1]))
7    bh = Normal(mu=tf.zeros(H), sigma=tf.ones(H))
8    by = Normal(mu=tf.zeros(1), sigma=tf.ones(1))
9
10   x = tf.placeholder(tf.float32, [None, D])
11   h = tf.scan(rnn_cell, x, initializer=tf.zeros(H))
12   y = Normal(mu=tf.matmul(h, Wy) + by, sigma=1.0)
```

## Composing Inferences

Core to Edward's design is that inference can be written as a collection of separate inference programs. Below we demonstrate variational EM.

```
1    qbeta = PointMass(params=tf.Variable(tf.zeros([K, D])))
2    qz = Categorical(logits=tf.Variable(tf.zeros[N, K]))
3
4    inference_e = ed.VariationalInference({z: qz}, data={x: x_data, beta: qbeta})
5    inference_m = ed.MAP({beta: qbeta}, data={x: x_data, z: qz})
6
7    for _ in range(10000):
8      inference_e.update()
9      inference_m.update()
```
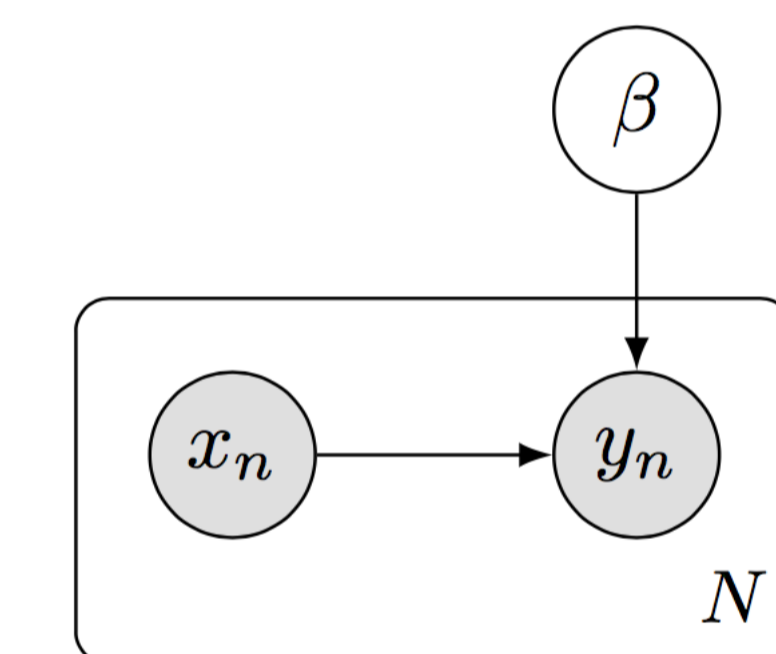
## Experiments: Recent Methods in Variational Inference

| Inference method | Negative log-likelihood |
|---|---|
| Variational auto-encoder (VAE) [2] | $\le 88.2$ |
| VAE without analytic KL | $\le 89.4$ |
| VAE with analytic entropy | $\le 88.1$ |
| VAE with score function gradient | $\le 87.9$ |
| Normalizing flows [4] | $\le 85.8$ |
| Hierarchical variational model [3] | $\le 85.4$ |
| Importance-weighted auto-encoders ($K = 50$) [1] | $\le 86.3$ |
| HVM with IWAE objective ($K = 5$) | $\le 85.2$ |
| Rényi divergence ($\alpha = -1$) | $\le 140.5$ |

Inference methods for a probabilistic decoder on binarized MNIST. The Edward PPL enables fast experimentation with many algorithms.

## Experiments: GPU-accelerated Hamiltonian Monte Carlo



```
1    # Model
2    x = tf.Variable(x_data, trainable=False)
3    beta = Normal(mu=tf.zeros(D), sigma=tf.ones(D))
4    y = Bernoulli(logits=tf.dot(x, beta))
5
6    # Inference
7    qbeta = Empirical(params=tf.Variable(tf.zeros([T, D])))
8    inference = ed.HMC({beta: qbeta}, data={y: y_data})
9    inference.run(step_size=0.5 / N, n_steps=100)
```

We apply Bayesian logistic regression to Covertype ($N = 581012$, $D = 54$). 12-core Intel i7-5930K CPU at 3.50GHz, a NVIDIA Titan X (Maxwell) GPU. We compare the runtime of HMC for 100 iterations (and same settings).

| Probabilistic programming system | Runtime (s) |
|---|---|
| Handwritten NumPy (1 CPU) | 534 |
| Stan (1 CPU) | 171 |
| PyMC3 (12 CPU) | 30.0 |
| **Edward (12 CPU)** | **8.2** |
| Handwritten TensorFlow (GPU) | 5.0 |
| **Edward (GPU)** | **4.9 (35x faster than Stan)** |

Edward (GPU) is significantly faster than other systems. In addition, Edward has no overhead: it is as fast as handwritten TensorFlow.

## References

[1]  Burda, Y., Grosse, R., and Salakhutdinov, R. (2016). Importance weighted autoencoders. In *International Conference on Learning Representations*.

[2]  Kingma, D. P. and Welling, M. (2014). Auto-encoding variational Bayes. In *International Conference on Learning Representations*.

[3]  Ranganath, R., Tran, D., and Blei, D. M. (2016). Hierarchical variational models. In *International Conference on Machine Learning*.

[4]  Rezende, D. J. and Mohamed, S. (2015). Variational inference with normalizing flows. In *International Conference on Machine Learning*.