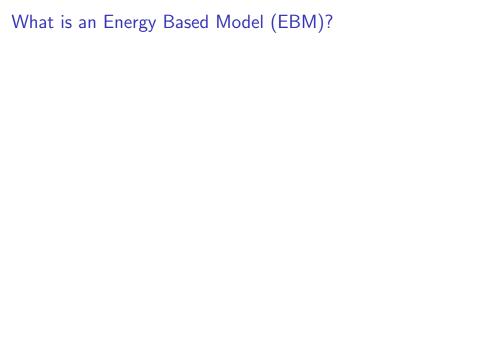
Your Classifier is Secretly an Energy Based Model and You Should Treat it Like One¹

January 30, 2020

¹Will Grathwohl et al., ICLR 2020



Any pdf $p_{\theta}(x)$ can be written down as:

$$p_{\theta}(\mathbf{x}) = \frac{\exp\left(-E_{\theta}(\mathbf{x})\right)}{Z_{\theta}} \tag{1}$$

Any pdf $p_{\theta}(x)$ can be written down as:

$$p_{\theta}(\mathbf{x}) = \frac{\exp\left(-E_{\theta}(\mathbf{x})\right)}{Z_{\theta}} \tag{1}$$

$$Z_{\theta} = \int_{\mathbf{x}} \exp\left(-E_{\theta}(\mathbf{x})\right) d\mathbf{x}$$

Any pdf $p_{\theta}(x)$ can be written down as:

$$p_{\theta}(\mathbf{x}) = \frac{\exp\left(-E_{\theta}(\mathbf{x})\right)}{Z_{\theta}} \tag{1}$$

where Z_{θ} is a normalizing constant:

$$Z_{\theta} = \int_{\mathbf{x}} \exp\left(-E_{\theta}(\mathbf{x})\right) d\mathbf{x}$$

▶ This parametrization allows us to forget about $p_{\theta}(x)$ and work with $E_{\theta}(x)$

Any pdf $p_{\theta}(x)$ can be written down as:

$$p_{\theta}(\mathbf{x}) = \frac{\exp\left(-E_{\theta}(\mathbf{x})\right)}{Z_{\theta}} \tag{1}$$

$$Z_{\theta} = \int_{\mathbf{x}} \exp\left(-E_{\theta}(\mathbf{x})\right) d\mathbf{x}$$

- ▶ This parametrization allows us to forget about $p_{\theta}(x)$ and work with $E_{\theta}(x)$
- ▶ Function $E_{\theta}(x)$ is an *energy function* (hence the name EBM)

Any pdf $p_{\theta}(x)$ can be written down as:

$$p_{\theta}(\mathbf{x}) = \frac{\exp\left(-E_{\theta}(\mathbf{x})\right)}{Z_{\theta}} \tag{1}$$

$$Z_{\theta} = \int_{\mathbf{x}} \exp\left(-E_{\theta}(\mathbf{x})\right) d\mathbf{x}$$

- ▶ This parametrization allows us to forget about $p_{\theta}(x)$ and work with $E_{\theta}(x)$
- ▶ Function $E_{\theta}(x)$ is an *energy function* (hence the name EBM)
- ▶ Energy functions are *easier* to work with: you do not have any restrictions on $E_{\theta}(x)$ ($p_{\theta}(x)$ is nonnegative and integrates to 1).

Any pdf $p_{\theta}(x)$ can be written down as:

$$p_{\theta}(\mathbf{x}) = \frac{\exp\left(-E_{\theta}(\mathbf{x})\right)}{Z_{\theta}} \tag{1}$$

$$Z_{\theta} = \int_{\mathbf{x}} \exp\left(-E_{\theta}(\mathbf{x})\right) d\mathbf{x}$$

- ▶ This parametrization allows us to forget about $p_{\theta}(x)$ and work with $E_{\theta}(x)$
- ▶ Function $E_{\theta}(x)$ is an *energy function* (hence the name EBM)
- ▶ Energy functions are *easier* to work with: you do not have any restrictions on $E_{\theta}(x)$ ($p_{\theta}(x)$ is nonnegative and integrates to 1).
- ▶ Energy functions are *harder* to work with: Z_{θ} is hard or impossible to compute, so we can't optimize EBM for θ

▶ We perform standard MLE: maximize $\log p_{\theta}(x_i)$ at our points x_i .

- ▶ We perform standard MLE: maximize $\log p_{\theta}(x_i)$ at our points x_i .
- ▶ By optimizing log $p_{\theta}(x)$ we also optimize $E_{\theta}(x)$.

- ▶ We perform standard MLE: maximize $\log p_{\theta}(x_i)$ at our points x_i .
- ▶ By optimizing log $p_{\theta}(x)$ we also optimize $E_{\theta}(x)$.
- ▶ We can't compute $\log p_{\theta}(x)$, but we can (approximately) compute its log-derivative via:

$$\frac{\partial \log p_{\theta}(\mathbf{x})}{\partial \theta} = \mathbb{E}_{p_{\theta}(\mathbf{x}')} \left[\frac{\partial E_{\theta}(\mathbf{x}')}{\partial \theta} \right] - \frac{\partial E_{\theta}(\mathbf{x})}{\partial \theta}$$
(2)

- ▶ We perform standard MLE: maximize $\log p_{\theta}(x_i)$ at our points x_i .
- ▶ By optimizing log $p_{\theta}(x)$ we also optimize $E_{\theta}(x)$.
- ▶ We can't compute $\log p_{\theta}(x)$, but we can (approximately) compute its log-derivative via:

$$\frac{\partial \log p_{\theta}(\mathbf{x})}{\partial \theta} = \mathbb{E}_{p_{\theta}(\mathbf{x}')} \left[\frac{\partial E_{\theta}(\mathbf{x}')}{\partial \theta} \right] - \frac{\partial E_{\theta}(\mathbf{x})}{\partial \theta}$$
(2)

Proof: follows directly from the definition.

- ▶ We perform standard MLE: maximize $\log p_{\theta}(x_i)$ at our points x_i .
- ▶ By optimizing log $p_{\theta}(x)$ we also optimize $E_{\theta}(x)$.
- We can't compute $\log p_{\theta}(x)$, but we can (approximately) compute its log-derivative via:

$$\frac{\partial \log p_{\theta}(\mathbf{x})}{\partial \theta} = \mathbb{E}_{p_{\theta}(\mathbf{x}')} \left[\frac{\partial E_{\theta}(\mathbf{x}')}{\partial \theta} \right] - \frac{\partial E_{\theta}(\mathbf{x})}{\partial \theta}$$
(2)

- Proof: follows directly from the definition.
- But a new problem arises: how to compute expectation?

- ▶ We perform standard MLE: maximize $\log p_{\theta}(x_i)$ at our points x_i .
- ▶ By optimizing log $p_{\theta}(x)$ we also optimize $E_{\theta}(x)$.
- ▶ We can't compute $\log p_{\theta}(x)$, but we can (approximately) compute its log-derivative via:

$$\frac{\partial \log p_{\theta}(\mathbf{x})}{\partial \theta} = \mathbb{E}_{p_{\theta}(\mathbf{x}')} \left[\frac{\partial E_{\theta}(\mathbf{x}')}{\partial \theta} \right] - \frac{\partial E_{\theta}(\mathbf{x})}{\partial \theta}$$
(2)

- Proof: follows directly from the definition.
- But a new problem arises: how to compute expectation?
- Answer: we can't do that directly, so let's use SGLD sampling.

Motivation:

▶ Imagine we have a "magic" function $p_{\theta}(x)$ that takes an input x and tells us its density.

- ▶ Imagine we have a "magic" function $p_{\theta}(x)$ that takes an input x and tells us its density.
- ▶ Unfortunately, we cannot sample new *x*-s from it directly, i.e. $p_{\theta}(x)$ does not support such operations.

- ▶ Imagine we have a "magic" function $p_{\theta}(x)$ that takes an input x and tells us its density.
- ▶ Unfortunately, we cannot sample new *x*-s from it directly, i.e. $p_{\theta}(x)$ does not support such operations.
- ▶ But we really want to sample from it! What should we do?

- ▶ Imagine we have a "magic" function $p_{\theta}(x)$ that takes an input x and tells us its density.
- ▶ Unfortunately, we cannot sample new *x*-s from it directly, i.e. $p_{\theta}(x)$ does not support such operations.
- ▶ But we really want to sample from it! What should we do?
- SGLD is a method to sample from such functions.

- ▶ Imagine we have a "magic" function $p_{\theta}(x)$ that takes an input x and tells us its density.
- ▶ Unfortunately, we cannot sample new *x*-s from it directly, i.e. $p_{\theta}(x)$ does not support such operations.
- ▶ But we really want to sample from it! What should we do?
- ▶ SGLD is a method to sample from such functions.
- Mhat is cool about SGLD is that it allows to use *unnormalized* density $\tilde{p}(x)$ (a very important property in our case!).

SGLD in a (very) few words:

1. Start from any random point x_0 .

- 1. Start from any random point x_0 .
- 2. Compute the gradient of $\nabla \log p_{\theta}(x_0)$.

- 1. Start from any random point x_0 .
- 2. Compute the gradient of $\nabla \log p_{\theta}(x_0)$.
- 3. Perform a small step in that direction.

- 1. Start from any random point x_0 .
- 2. Compute the gradient of $\nabla \log p_{\theta}(x_0)$.
- 3. Perform a small step in that direction.
- 4. Add a little bit of noise.

- 1. Start from any random point x_0 .
- 2. Compute the gradient of $\nabla \log p_{\theta}(x_0)$.
- 3. Perform a small step in that direction.
- 4. Add a little bit of noise.
- 5. Obtain x_{t+1} . Accept it if:

$$\frac{p(x_t|x_{t+1})p(x_t)}{p(x_{t+1}|x_t)p(x_{t+1})} < u, u \sim \mathcal{U}[0, 1]$$
(3)

SGLD in a (very) few words:

- 1. Start from any random point x_0 .
- 2. Compute the gradient of $\nabla \log p_{\theta}(x_0)$.
- 3. Perform a small step in that direction.
- 4. Add a little bit of noise.
- 5. Obtain x_{t+1} . Accept it if:

$$\frac{p(x_t|x_{t+1})p(x_t)}{p(x_{t+1}|x_t)p(x_{t+1})} < u, u \sim \mathcal{U}[0,1]$$
(3)

In other words, run the following recurrence a lot of times:

$$\mathbf{x}_0 \sim p_0(\mathbf{x}), \quad \mathbf{x}_{i+1} = \mathbf{x}_i - \frac{\alpha}{2} \frac{\partial E_{\theta}(\mathbf{x}_i)}{\partial \mathbf{x}_i} + \epsilon, \quad \epsilon \sim \mathcal{N}(0, \alpha)$$
 (4)

SGLD in a (very) few words:

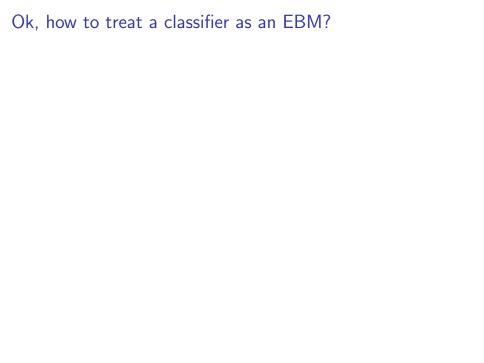
- 1. Start from any random point x_0 .
- 2. Compute the gradient of $\nabla \log p_{\theta}(x_0)$.
- 3. Perform a small step in that direction.
- 4. Add a little bit of noise.
- 5. Obtain x_{t+1} . Accept it if:

$$\frac{p(x_t|x_{t+1})p(x_t)}{p(x_{t+1}|x_t)p(x_{t+1})} < u, u \sim \mathcal{U}[0,1]$$
(3)

In other words, run the following recurrence a lot of times:

$$\mathbf{x}_0 \sim p_0(\mathbf{x}), \quad \mathbf{x}_{i+1} = \mathbf{x}_i - \frac{\alpha}{2} \frac{\partial E_{\theta}(\mathbf{x}_i)}{\partial \mathbf{x}_i} + \epsilon, \quad \epsilon \sim \mathcal{N}(0, \alpha)$$
 (4)

Justification: if your model and your choice of α are "good" (in some sense), then after infinite amount of steps your samples will be "good" (in some sense).



▶ Remember that any function can be used as $E_{\theta}(x)$

- ▶ Remember that *any* function can be used as $E_{\theta}(x)$
- ▶ Let $f_{\theta}(x)$ be our neural network.

- ▶ Remember that any function can be used as $E_{\theta}(x)$
- ▶ Let $f_{\theta}(x)$ be our neural network.
- ▶ Let $f_{\theta}(x)[y]$ be the y-th logit (i.e. logit value for class y).

- ▶ Remember that any function can be used as $E_{\theta}(x)$
- ▶ Let $f_{\theta}(x)$ be our neural network.
- ▶ Let $f_{\theta}(x)[y]$ be the *y*-th logit (i.e. logit value for class *y*).
- ▶ Let our examples be s = (x, y)

- ▶ Remember that any function can be used as $E_{\theta}(x)$
- ▶ Let $f_{\theta}(x)$ be our neural network.
- ▶ Let $f_{\theta}(x)[y]$ be the *y*-th logit (i.e. logit value for class *y*).
- ▶ Let our examples be s = (x, y)
- We define energy $E_{\theta}(s) = E_{\theta}(x, y) = -f_{\theta}(x)[y]$.

- ▶ Remember that any function can be used as $E_{\theta}(x)$
- ▶ Let $f_{\theta}(x)$ be our neural network.
- Let $f_{\theta}(x)[y]$ be the y-th logit (i.e. logit value for class y).
- ▶ Let our examples be s = (x, y)
- ▶ We define energy $E_{\theta}(s) = E_{\theta}(x, y) = -f_{\theta}(x)[y]$.
- ▶ Using $E_{\theta}(x,y)$ we can define $p_{\theta}(x)$:

$$p_{\theta}(\mathbf{x}) = \sum_{y} p_{\theta}(\mathbf{x}, y) = \frac{1}{Z_{\theta}} \sum_{y} \exp(f_{\theta}(\mathbf{x})[y])$$

(5)

Ok, how to treat a classifier as an EBM?

- ▶ Remember that any function can be used as $E_{\theta}(x)$
- ▶ Let $f_{\theta}(x)$ be our neural network.
- ▶ Let $f_{\theta}(x)[y]$ be the y-th logit (i.e. logit value for class y).
- ▶ Let our examples be s = (x, y)
- ▶ We define energy $E_{\theta}(s) = E_{\theta}(x, y) = -f_{\theta}(x)[y]$.
- Using $E_{\theta}(x,y)$ we can define $p_{\theta}(x)$:

$$p_{\theta}(\mathbf{x}) = \sum_{y} p_{\theta}(\mathbf{x}, y) = \frac{1}{Z_{\theta}} \sum_{y} \exp(f_{\theta}(\mathbf{x})[y])$$
 (5)

• Using $p_{\theta}(x)$ we can define $E_{\theta}(x)$:

$$E_{\theta}(\mathbf{x}) = -\log \sum_{\mathbf{x}} \exp\left(f_{\theta}(\mathbf{x})[y]\right) \tag{6}$$

Ok, how to treat a classifier as an EBM?

- ▶ Remember that any function can be used as $E_{\theta}(x)$
- ▶ Let $f_{\theta}(x)$ be our neural network.
- ▶ Let $f_{\theta}(x)[y]$ be the *y*-th logit (i.e. logit value for class *y*).
- ▶ Let our examples be s = (x, y)
- ▶ We define energy $E_{\theta}(s) = E_{\theta}(x, y) = -f_{\theta}(x)[y]$.
- ▶ Using $E_{\theta}(x,y)$ we can define $p_{\theta}(x)$:

$$p_{\theta}(\mathbf{x}) = \sum_{y} p_{\theta}(\mathbf{x}, y) = \frac{1}{Z_{\theta}} \sum_{y} \exp\left(f_{\theta}(\mathbf{x})[y]\right)$$
 (5)

• Using $p_{\theta}(x)$ we can define $E_{\theta}(x)$:

$$E_{\theta}(\mathbf{x}) = -\log \sum_{\mathbf{x}} \exp(f_{\theta}(\mathbf{x})[y])$$
 (6)

▶ Higher the logits — lower the energy of x.

Ok, how to treat a classifier as an EBM?

- ▶ Remember that any function can be used as $E_{\theta}(x)$
- ▶ Let $f_{\theta}(x)$ be our neural network.
- Let $f_{\theta}(x)[y]$ be the y-th logit (i.e. logit value for class y).
- ▶ Let our examples be s = (x, y)
- ▶ We define energy $E_{\theta}(s) = E_{\theta}(x, y) = -f_{\theta}(x)[y]$.
- ▶ Using $E_{\theta}(x, y)$ we can define $p_{\theta}(x)$:

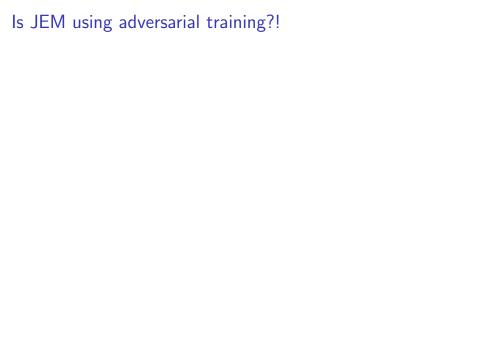
$$p_{\theta}(\mathbf{x}) = \sum_{y} p_{\theta}(\mathbf{x}, y) = \frac{1}{Z_{\theta}} \sum_{y} \exp\left(f_{\theta}(\mathbf{x})[y]\right)$$
 (5)

• Using $p_{\theta}(x)$ we can define $E_{\theta}(x)$:

$$E_{\theta}(\mathbf{x}) = -\log \sum_{y} \exp \left(f_{\theta}(\mathbf{x})[y] \right)$$
 (6)

- ▶ Higher the logits lower the energy of *x*.
- Our final loss looks like:

$$\mathcal{L}(x) = \log p(x, y) = \underbrace{\log p(y|x)}_{\text{classification loss}} + \underbrace{\log p(x)}_{\text{generative loss}}$$
(7)



Let's take a closer look to the gradient of the part of the loss:

$$\frac{\partial \log p_{\theta}(\mathbf{x})}{\partial \theta} = \mathbb{E}_{p_{\theta}(\mathbf{x}')} \left[\frac{\partial E_{\theta}(\mathbf{x}')}{\partial \theta} \right] - \frac{\partial E_{\theta}(\mathbf{x})}{\partial \theta}$$
(8)

Let's take a closer look to the gradient of the part of the loss:

$$\frac{\partial \log p_{\theta}(\mathbf{x})}{\partial \theta} = \mathbb{E}_{p_{\theta}(\mathbf{x}')} \left[\frac{\partial E_{\theta}(\mathbf{x}')}{\partial \theta} \right] - \frac{\partial E_{\theta}(\mathbf{x})}{\partial \theta}$$
(8)

▶ If we use a single sample to approximate expectation and compute the gradient across the dataset our gradient we'll look like:

$$\frac{\partial \log p_{\theta}(\mathbf{X})}{\partial \theta} \approx \frac{1}{n} \sum_{i=1}^{n} \left(\frac{\partial E_{\theta} \left(\mathbf{x}^{\mathsf{fake}} \right)}{\partial \theta} - \frac{\partial E_{\theta} \left(\mathbf{x}^{\mathsf{real}} \right)}{\partial \theta} \right) \tag{9}$$

Let's take a closer look to the gradient of the part of the loss:

$$\frac{\partial \log p_{\theta}(\mathbf{x})}{\partial \theta} = \mathbb{E}_{p_{\theta}(\mathbf{x}')} \left[\frac{\partial E_{\theta}(\mathbf{x}')}{\partial \theta} \right] - \frac{\partial E_{\theta}(\mathbf{x})}{\partial \theta}$$
(8)

▶ If we use a single sample to approximate expectation and compute the gradient across the dataset our gradient we'll look like:

$$\frac{\partial \log p_{\theta}(\mathbf{X})}{\partial \theta} \approx \frac{1}{n} \sum_{i=1}^{n} \left(\frac{\partial E_{\theta} \left(\mathbf{x}^{\mathsf{fake}} \right)}{\partial \theta} - \frac{\partial E_{\theta} \left(\mathbf{x}^{\mathsf{real}} \right)}{\partial \theta} \right) \tag{9}$$

And this gradient looks very similar to the gradient of the following loss:

$$\mathcal{L}(\theta) = \underset{p_{\theta}(x)}{\mathbb{E}} [E_{\theta}(x)] - \underset{p_{\text{data}}(x)}{\mathbb{E}} [E_{\theta}(x)]$$
 (10)

Let's take a closer look to the gradient of the part of the loss:

$$\frac{\partial \log p_{\theta}(\mathbf{x})}{\partial \theta} = \mathbb{E}_{p_{\theta}(\mathbf{x}')} \left[\frac{\partial E_{\theta}(\mathbf{x}')}{\partial \theta} \right] - \frac{\partial E_{\theta}(\mathbf{x})}{\partial \theta}$$
(8)

▶ If we use a single sample to approximate expectation and compute the gradient across the dataset our gradient we'll look like:

$$\frac{\partial \log p_{\theta}(\mathbf{X})}{\partial \theta} \approx \frac{1}{n} \sum_{i=1}^{n} \left(\frac{\partial E_{\theta} \left(\mathbf{x}^{\mathsf{fake}} \right)}{\partial \theta} - \frac{\partial E_{\theta} \left(\mathbf{x}^{\mathsf{real}} \right)}{\partial \theta} \right) \tag{9}$$

► And this gradient looks very similar to the gradient of the following loss:

$$\mathcal{L}(\theta) = \underset{p_{\theta}(x)}{\mathbb{E}} \left[E_{\theta}(x) \right] - \underset{p_{\text{data}}(x)}{\mathbb{E}} \left[E_{\theta}(x) \right] \tag{10}$$

Which is a discriminator loss!

Let's take a closer look to the gradient of the part of the loss:

$$\frac{\partial \log p_{\theta}(\mathbf{x})}{\partial \theta} = \mathbb{E}_{p_{\theta}(\mathbf{x}')} \left[\frac{\partial E_{\theta}(\mathbf{x}')}{\partial \theta} \right] - \frac{\partial E_{\theta}(\mathbf{x})}{\partial \theta}$$
(8)

▶ If we use a single sample to approximate expectation and compute the gradient across the dataset our gradient we'll look like:

$$\frac{\partial \log p_{\theta}(\mathbf{X})}{\partial \theta} \approx \frac{1}{n} \sum_{i=1}^{n} \left(\frac{\partial E_{\theta} \left(\mathbf{x}^{\mathsf{fake}} \right)}{\partial \theta} - \frac{\partial E_{\theta} \left(\mathbf{x}^{\mathsf{real}} \right)}{\partial \theta} \right) \tag{9}$$

And this gradient looks very similar to the gradient of the following loss:

$$\mathcal{L}(\theta) = \underset{p_{\theta}(x)}{\mathbb{E}} \left[E_{\theta}(x) \right] - \underset{p_{\text{data}}(x)}{\mathbb{E}} \left[E_{\theta}(x) \right] \tag{10}$$

- Which is a discriminator loss!
- ▶ Reformulation of GANs via EBMs is not new.

Let's take a closer look to the gradient of the part of the loss:

$$\frac{\partial \log p_{\theta}(\mathbf{x})}{\partial \theta} = \mathbb{E}_{p_{\theta}(\mathbf{x}')} \left[\frac{\partial E_{\theta}(\mathbf{x}')}{\partial \theta} \right] - \frac{\partial E_{\theta}(\mathbf{x})}{\partial \theta}$$
(8)

▶ If we use a single sample to approximate expectation and compute the gradient across the dataset our gradient we'll look like:

$$\frac{\partial \log p_{\theta}(\mathbf{X})}{\partial \theta} \approx \frac{1}{n} \sum_{i=1}^{n} \left(\frac{\partial E_{\theta} \left(\mathbf{x}^{\mathsf{fake}} \right)}{\partial \theta} - \frac{\partial E_{\theta} \left(\mathbf{x}^{\mathsf{real}} \right)}{\partial \theta} \right) \tag{9}$$

And this gradient looks very similar to the gradient of the following loss:

$$\mathcal{L}(\theta) = \underset{p_{\theta}(x)}{\mathbb{E}} \left[E_{\theta}(x) \right] - \underset{p_{\text{data}}(x)}{\mathbb{E}} \left[E_{\theta}(x) \right] \tag{10}$$

- Which is a discriminator loss!
- Reformulation of GANs via EBMs is not new.
- ► The key novelty of the paper is that we can do that "inside" the classifier.

▶ Disclaimer: results are really good.

- ▶ Disclaimer: results are really good.
- ▶ Authors built a model on SVHN/CIFAR10/CIFAR100 which:

- ▶ Disclaimer: results are really good.
- ▶ Authors built a model on SVHN/CIFAR10/CIFAR100 which:
 - has good classification accuracy

- Disclaimer: results are really good.
- ▶ Authors built a model on SVHN/CIFAR10/CIFAR100 which:
 - has good classification accuracy
 - ▶ has good IS/FID scores (i.e. has good samples)

- Disclaimer: results are really good.
- ▶ Authors built a model on SVHN/CIFAR10/CIFAR100 which:
 - has good classification accuracy
 - has good IS/FID scores (i.e. has good samples)
 - has much better calibrated predictions (i.e. is not overconfident unlike traditional classifiers)

- Disclaimer: results are really good.
- ▶ Authors built a model on SVHN/CIFAR10/CIFAR100 which:
 - has good classification accuracy
 - has good IS/FID scores (i.e. has good samples)
 - ► has much better calibrated predictions (i.e. is not overconfident unlike traditional classifiers)
 - has good out-of-distribution predictions (i.e. detecting when we feed OOD-samples into it)

- Disclaimer: results are really good.
- ▶ Authors built a model on SVHN/CIFAR10/CIFAR100 which:
 - has good classification accuracy
 - has good IS/FID scores (i.e. has good samples)
 - ► has much better calibrated predictions (i.e. is not overconfident unlike traditional classifiers)
 - has good out-of-distribution predictions (i.e. detecting when we feed OOD-samples into it)
 - is robust to adversarial attacks

- Disclaimer: results are really good.
- ▶ Authors built a model on SVHN/CIFAR10/CIFAR100 which:
 - has good classification accuracy
 - has good IS/FID scores (i.e. has good samples)
 - has much better calibrated predictions (i.e. is not overconfident unlike traditional classifiers)
 - has good out-of-distribution predictions (i.e. detecting when we feed OOD-samples into it)
 - is robust to adversarial attacks
 - ...and all of these at the same time!

- Disclaimer: results are really good.
- ▶ Authors built a model on SVHN/CIFAR10/CIFAR100 which:
 - has good classification accuracy
 - has good IS/FID scores (i.e. has good samples)
 - has much better calibrated predictions (i.e. is not overconfident unlike traditional classifiers)
 - has good out-of-distribution predictions (i.e. detecting when we feed OOD-samples into it)
 - is robust to adversarial attacks
 - ...and all of these at the same time!
- A huge limitation: SGLD makes both the training unstable and the model hard to scale.

- Disclaimer: results are really good.
- ▶ Authors built a model on SVHN/CIFAR10/CIFAR100 which:
 - has good classification accuracy
 - has good IS/FID scores (i.e. has good samples)
 - has much better calibrated predictions (i.e. is not overconfident unlike traditional classifiers)
 - has good out-of-distribution predictions (i.e. detecting when we feed OOD-samples into it)
 - is robust to adversarial attacks
 - ...and all of these at the same time!
- ▶ A huge limitation: SGLD makes both the training unstable and the model hard to scale.
- ► An approach was coined Joint Energy based Model (JEM) since we model (x, y) simultaneously.