

Your classifier is secretly an energy based  
model and you should treat it like one  
(ICLR 2020)

**2020-04-24**

**hyounguk.shon@kaist.ac.kr**

# Abstract

- This paper proposes Joint Energy-based Model. (JEM)
- JEM models the joint distribution  $p(\mathbf{x}, y)$  by reinterpreting softmax-based classifier architecture as an energy-based model (EBM).
- JEM is trained to estimate (unnormalized)  $p(\mathbf{x})$  and  $p(y|\mathbf{x})$ . Therefore, JEM can be trained on both unlabeled and labeled data.
- JEM improves classifiers in terms of calibration, robustness, OOD detection, while enabling it to generate samples.
- JEM achieves performance rivaling the SOTA in both generative and discriminative tasks with a single model.

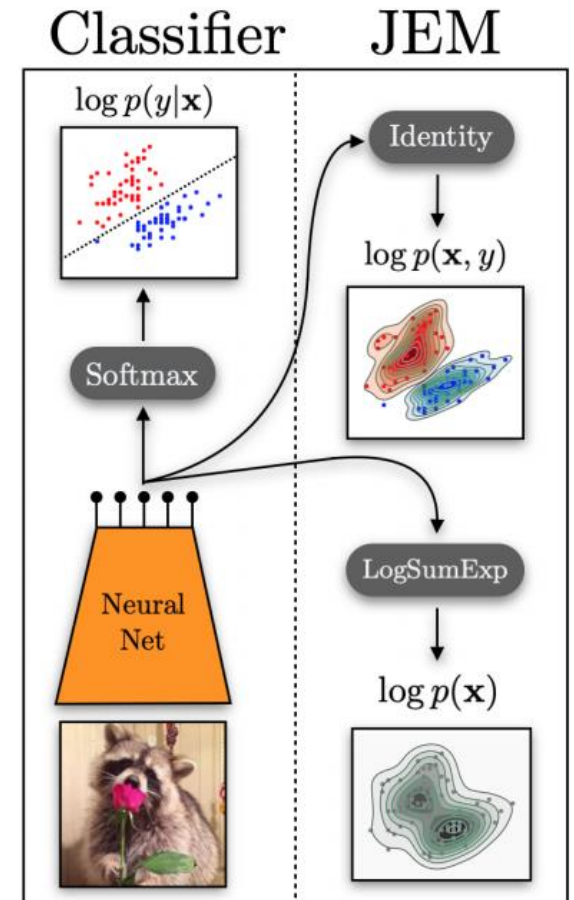


Figure 1: Visualization of our method, JEM, which defines a joint EBM from classifier architectures.

# Introduction to EBM (I)

- Energy-based Model (EBM) [1]:
  - EBM estimates a probability density function by introducing an energy function  $E_\theta(\mathbf{x})$ .
  - Any arbitrary  $p(\mathbf{x})$  can be modeled by a function  $E_\theta(\mathbf{x}) : \mathbb{R}^D \rightarrow \mathbb{R}$ .

$$p_\theta(\mathbf{x}) \equiv \frac{\exp(-E_\theta(\mathbf{x}))}{Z_\theta} \quad \text{where } Z_\theta = \int_{\mathbf{x}} \exp(-E_\theta(\mathbf{x}))$$

- Inference and Training on EBM:
  - Notice that  $Z_\theta$  is a normalizing constant. For most choices of  $E_\theta(\mathbf{x})$ , it is difficult to compute or estimate  $Z_\theta$ . Therefore it is intractable to compute the normalized likelihood.
  - Therefore, we must rely on other methods to train EBMs.

[1] Y LeCun et al., A tutorial on energy-based learning. Predicting structured data, 2006.

# Introduction to EBM (II)

- The gradient of the log-likelihood is expressed as follows:

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

- It is hard to compute the expectation because we cannot easily draw samples from  $p_{\theta}(\mathbf{x})$ . So we must resort to MCMC sampling to estimate the expectation.
- Recent successes on training large-scale EBM are based on Stochastic Gradient Langevin Dynamics (SGLD) sampler [2], which draws a sample through iterative refinement:

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

[2] M Welling et al., Bayesian learning via stochastic gradient langevin dynamics. ICML, 2011.

# What your classifier is hiding (I)

- Standard softmax architectures model conditional probability by predicting logits of likelihood and passing through a softmax layer which normalizes the likelihood.

$$p(y|\mathbf{x}) = \frac{\exp(f_{\theta}(\mathbf{x})[y])}{\sum_y \exp(f_{\theta}(\mathbf{x})[y])}$$

- $f_{\theta}(\mathbf{x})[y]$  indicates the logit value of the input  $\mathbf{x}$  that corresponds to a class  $y$ .

# What your classifier is hiding (II)

- Our goal is to model the joint probability  $p(\mathbf{x}, y)$  in the form of an EBM.

$$p_{\theta}(\mathbf{x}, y) = \frac{\exp(-E_{\theta}(\mathbf{x}, y))}{Z_{\theta}} \quad \text{where} \quad \int_{\mathbf{x}, y} \exp(-E_{\theta}(\mathbf{x}, y))$$

- First we factorize the joint probability and replace  $p(\mathbf{x})$  by marginalizing:

$$p(\mathbf{x}, y) = p(\mathbf{x})p(y|\mathbf{x}) = \left( \sum_y p(\mathbf{x}, y) \right) p(y|\mathbf{x})$$

- If we choose  $E_{\theta}(\mathbf{x}, y) = -f_{\theta}(\mathbf{x})[y]$  and plug into EBM, we get:

$$\frac{\exp(f_{\theta}(\mathbf{x})[y])}{Z_{\theta}} = \frac{\sum_y \exp(f_{\theta}(\mathbf{x})[y])}{Z_{\theta}} p(y|\mathbf{x})$$

- This leads us back to the softmax classifier:

$$p(y|\mathbf{x}) = \frac{\exp(f_{\theta}(\mathbf{x})[y])}{\sum_y \exp(f_{\theta}(\mathbf{x})[y])}$$

- In other words, we have revealed a generative model hidden within every standard discriminative model!

# What your classifier is hiding (III)

- Notice that in standard softmax classifiers, shifting the logits does not affect the prediction. However with JEM, shifting the logits will affect  $p_\theta(\mathbf{x}, y)$ . This is because JEM utilizes the extra degree of freedom hidden in softmax classifiers to model  $p_\theta(\mathbf{x})$ .

$$\text{softmax}(\mathbf{h}) = \frac{\exp(h_i)}{\sum_y \exp(h_i)} = \frac{\exp(h_i + \alpha)}{\sum_y \exp(h_i + \alpha)}$$

- Overview of the results:

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

$$\log p_\theta(\mathbf{x}, y) = f_\theta(\mathbf{x})[y] - \log Z_\theta$$

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

$$\log p_\theta(\mathbf{x}) = \text{LogSumExp}_y f_\theta(\mathbf{x})[y] - \log Z_\theta$$

$$p_\theta(y|\mathbf{x}) = \frac{\exp(f_\theta(\mathbf{x})[y])}{\sum_y \exp(f_\theta(\mathbf{x})[y])}$$

- Where,

$$E_\theta(\mathbf{x}, y) = -f_\theta(\mathbf{x})[y]$$

$$E_\theta(\mathbf{x}) = -\text{LogSumExp}_y f_\theta(\mathbf{x})[y]$$

$$Z_\theta = \int_{\mathbf{x}, y} \exp(f_\theta(\mathbf{x})[y])$$

# Optimization (I)

- We optimize  $p_\theta(y|\mathbf{x})$  using cross-entropy, and optimize  $\log p_\theta(\mathbf{x})$  using SGLD to estimate the expectation.

$$\mathcal{L}_{\text{classifier}}(\mathbf{x}, y) = \text{CrossEntropy}(f_\theta(\mathbf{x}), y)$$

$$\mathcal{L}_{\text{generator}}(\mathbf{x}) = -\log p_\theta(\mathbf{x})$$

$$\text{objective} = \mathcal{L}_{\text{classifier}}(\mathbf{x}, y) + \mathcal{L}_{\text{generator}}(\mathbf{x})$$

- Following [3] we use persistent contrastive divergence [4] to estimate the expectation since it gives an order of magnitude savings in computation compared to seeding new chains at each iteration as in Nijakamp et al., this cost comes at the cost of decreased training stability. The trade-offs are discussed in appendix H.2.

[3] Y Du et al., Implicit generation and generalization in energy-based models. arXiv, 2019.

[4] T Tieleman. Training restricted boltzmann machines using approximations to the likelihood gradient. ICML, 2008



# Optimization (II)

---

**Algorithm 1** JEM training: Given network  $f_\theta$ , SGLD step-size  $\alpha$ , SGLD noise  $\sigma$ , replay buffer  $B$ , SGLD steps  $\eta$ , reinitialization frequency  $\rho$

---

```
1: while not converged do
2:   Sample  $\mathbf{x}$  and  $y$  from dataset
3:    $L_{\text{clf}}(\theta) = \text{xent}(f_\theta(\mathbf{x}), y)$ 
4:   Sample  $\hat{\mathbf{x}}_0 \sim B$  with probability  $1 - \rho$ , else  $\hat{\mathbf{x}}_0 \sim \mathcal{U}(-1, 1)$  ▷ Initialize SGLD
5:   for  $t \in [1, 2, \dots, \eta]$  do ▷ SGLD
6:      $\hat{\mathbf{x}}_t = \hat{\mathbf{x}}_{t-1} + \alpha \cdot \frac{\partial \text{LogSumExp}_{y'}(f_\theta(\hat{\mathbf{x}}_{t-1})[y'])}{\partial \hat{\mathbf{x}}_{t-1}} + \sigma \cdot \mathcal{N}(0, I)$ 
7:   end for
8:    $L_{\text{gen}}(\theta) = \text{LogSumExp}_{y'}(f(\mathbf{x})[y']) - \text{LogSumExp}_{y'}(f(\hat{\mathbf{x}}_t)[y'])$  ▷ Surrogate for Eq 2
9:    $L(\theta) = L_{\text{clf}}(\theta) + L_{\text{gen}}(\theta)$ 
10:  Obtain gradients  $\frac{\partial L(\theta)}{\partial \theta}$  for training
11:  Add  $\hat{\mathbf{x}}_t$  to  $B$ 
12: end while
```

---

# Applications

- Achieved performance rivaling to SOTA in both discriminative and generative modeling.
- Added benefits to discriminative tasks: uncertainty qualification, OOD detection, robustness to adversarial examples.
- All architectures are based on WideResNet with BN removed to ensure deterministic output. This slightly increases error of a WRN-28-10 from 4.2% to 6.4% on CIFAR-10 and from 2.3% to 3.4% on SVHN.

# Hybrid Modeling (I)

- A classifier architecture can be trained as an EBM to show competitive performance as both a classifier and generative model.
- JEM can achieve near-SOTA in both classification and generation, while outperforming other hybrid models. (CIFAR-10)

Class	Model	Accuracy% $\uparrow$	IS $\uparrow$	FID $\downarrow$
<b>Hybrid</b>	Residual Flow	70.3	3.6	46.4
	Glow	67.6	3.92	48.9
	IGEBM	49.1	8.3	<b>37.9</b>
	JEM $p(x y)$ factored	30.1	6.36	61.8
	JEM (Ours)	<b>92.9</b>	<b>8.76</b>	38.4
<b>Disc.</b>	Wide-Resnet	95.8	N/A	N/A
<b>Gen.</b>	SNGAN	N/A	8.59	25.5
	NCSN	N/A	8.91	25.32

Table 1: CIFAR10 Hybrid modeling Results. Residual Flow (Chen et al., 2019), Glow (Kingma & Dhariwal, 2018), IGEBM (Du & Mordatch, 2019), SNGAN (Miyato et al., 2018), NCSN (Song & Ermon, 2019)

# Hybrid Modeling (II)

- Class-conditional samples generated from JEM:



- The paper includes further interesting qualitative analyses in the appendix.



Figure 7: Class-conditional Samples. Left to right: CIFAR10, SVHN.

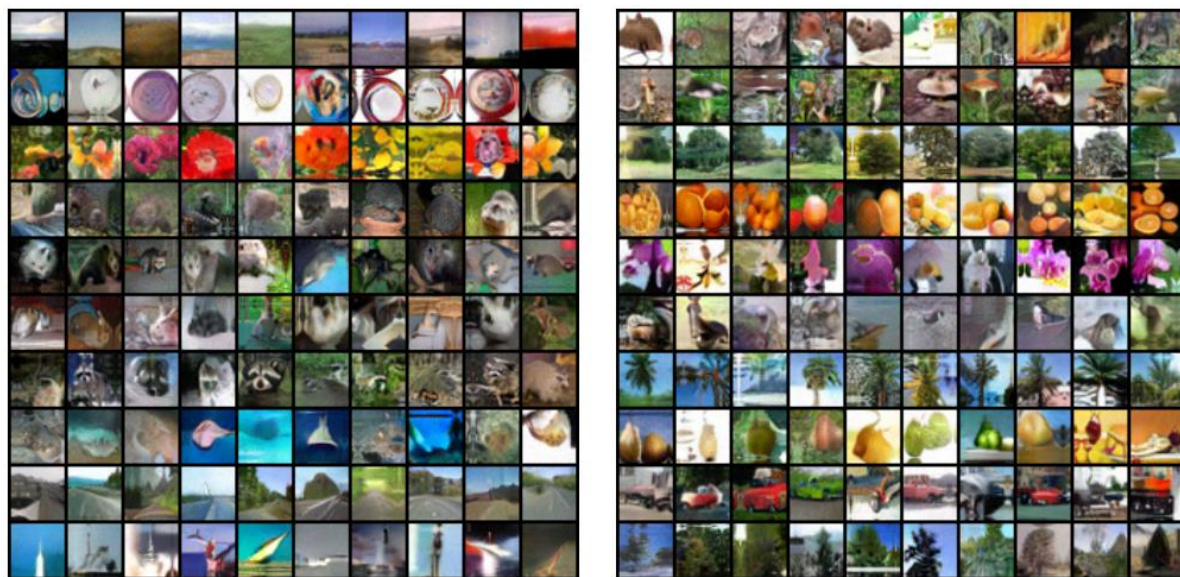
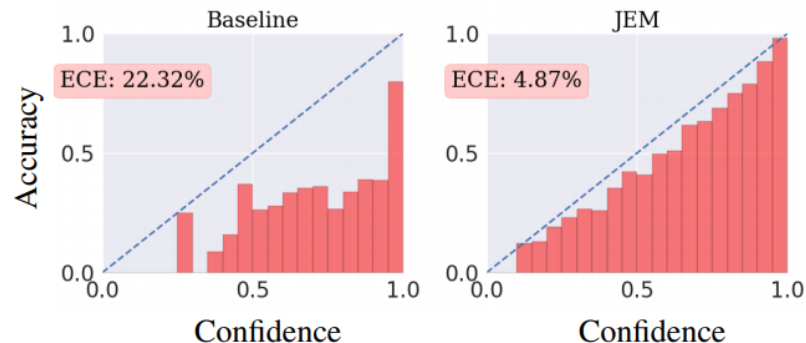


Figure 8: CIFAR100 Class-conditional Samples.



# Calibration

- A classifier is calibrated if its confidence  $\max_y p(y|\mathbf{x})$  is consistent with its error rate.
- Calibration is important in mission-critical applications, for example, where the confidence could be used to decide whether to defer the task to a human. Therefore, a well-calibrated but less accurate classifier can be considerably more useful than the other.
- Recent classifiers have grown more accurate but considerably less calibrated. [5] On the contrary, JEM notably improves classification while retaining high accuracy.



Model	Accuracy %	ECE %
Baseline	74.2	22.32
JEM	72.2	4.87

Figure 4: CIFAR100 calibration results. ECE = Expected Calibration Error (Guo et al., 2017), see Appendix E.1.

[5] C Guo et al., On calibration of modern neural networks. In ICML, 2017.

# Out-of-distribution Detection (I)

- In OOD detection, the goal to estimate an in-distributionness score  $s_{\theta}(\mathbf{x}) \in \mathbb{R}$ . We apply JEM to: input density estimation, predictive distribution estimation, approximate mass.
- Input density: this approach fits a density model  $p(\mathbf{x})$  on the data, and consider examples with low likelihood to be OOD. This approach is currently not competitive on high-dimensional data.
- [6] showed that tractable deep generative models [7, 8] can assign higher densities to OOD examples than in-distribution examples. Conversely, Du & Mordatch have shown that the likelihoods from EBMs can be reliably used as a predictor for OOD inputs. JEM consistently assigns higher likelihood to ID samples than OOD samples. One possible explanation for JEM's improvement over IGEBM is its ability to incorporate labeled information during training.

[6] E Nalisnick et al., Do deep generative models know what they don't know? arXiv, 2018.

[7] D P Kingma et al., Glow: Generative flow with invertible 1x1 convolutions. In NIPS, 2018.

[8] T Salimans et al., Pixelcnn++: Improving the pixelcnn with discretized logistic mixture likelihood and other modifications. arXiv, 2017.

# Out-of-distribution Detection (II)

- Predictive Distribution:  $s_{\theta}(\mathbf{x}) = \max_y p_{\theta}(y|\mathbf{x})$
- OOD detection by utilizing a classifier's predictive distribution have shown many success in [9, 10, 11, 12]. Since JEM is a competitive classifier, we find it performs on par or beyond the performance of a strong baseline classifier.

- [9] Y Gal et al., Dropout as a Bayesian approximation: Representing model uncertainty in deep learning. ICML, 2016.
- [10] K Wang et al., Adversarial distillation of bayesian neural network posteriors. ICML, 2018.
- [11] S Liang et al., Enhancing the reliability of out-of-distribution image detection in neural networks. arXiv 2017.
- [12] D Hendrycks et al., A baseline for detecting misclassified and out-of-distribution examples in neural networks. arXiv 2016.



# Out-of-distribution Detection (III)

- Approximate Mass:

$$s_{\theta}(\mathbf{x}) = - \left\| \frac{\partial \log p_{\theta}(\mathbf{x})}{\partial \mathbf{x}} \right\|_2$$

- It is possible that a sample has high likelihood  $p_{\theta}(\mathbf{x})$  under a distribution yet be nearly impossible to be sampled. Real samples lie in typical set, which is the area of high probability mass.
- For a high-likelihood sample outside of the typical set, we expect the density to change rapidly around it, thus the norm of the gradient of the log-density will be large compared to examples in the typical set. We propose an alternative OOD score based on this property:
- For EBM (JEM and IGEBM), we find this predictor greatly outperforms our own and other generative model's likelihoods. For tractable likelihood methods we find this predictor is anti-correlated with the model's likelihood and neither is reliable for OOD detection.

# Out-of-distribution Detection (IV)

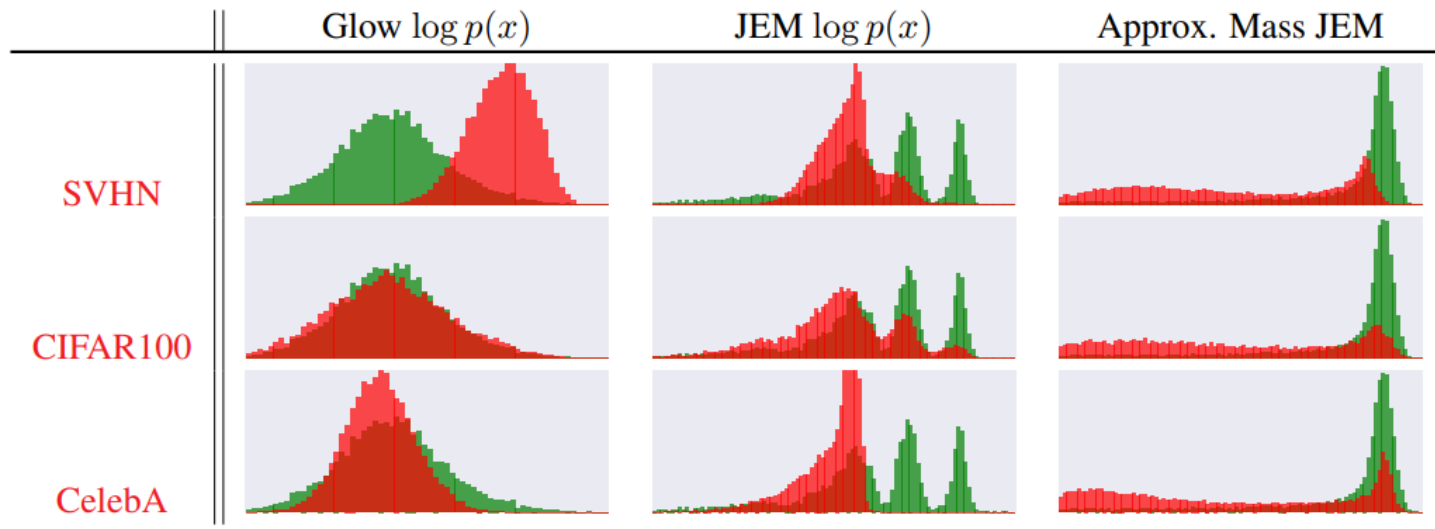


Table 2: Histograms for OOD detection. All models trained on **CIFAR10**. Green corresponds to the score on (in-distribution) **CIFAR10**, and red corresponds to the score on the OOD dataset.

In-distribution	Out-of-distribution
CIFAR10	SVHN
	CIFAR100
	CelebA

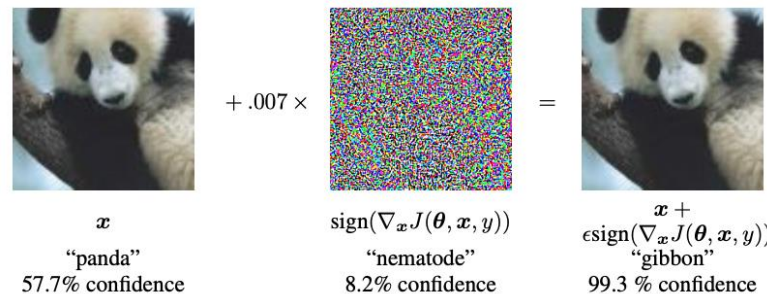
# Out-of-distribution Detection (V)

$s_{\theta}(\mathbf{x})$	Model	SVHN	CIFAR10 Interp	CIFAR100	CelebA
$\log p(\mathbf{x})$	Unconditional Glow	.05	.51	.55	.57
	Class-Conditional Glow	.07	.45	.51	.53
	IGEBM	.63	<b>.70</b>	.50	.70
	JEM (Ours)	<b>.67</b>	.65	<b>.67</b>	<b>.75</b>
$\max_y p(y \mathbf{x})$	Wide-ResNet	<b>.93</b>	<b>.77</b>	.85	.62
	Class-Conditional Glow	.64	.61	.65	.54
	IGEBM	.43	.69	.54	.69
	JEM (Ours)	.89	.75	<b>.87</b>	<b>.79</b>
$\left\  \frac{\partial \log p(\mathbf{x})}{\partial \mathbf{x}} \right\ $	Unconditional Glow	<b>.95</b>	.27	.46	.29
	Class-Conditional Glow	.47	.01	.52	.59
	IGEBM	.84	.65	.55	.66
	JEM (Ours)	.83	<b>.78</b>	<b>.82</b>	<b>.79</b>

Table 3: OOD Detection Results. Models trained on CIFAR10. Values are AUROC.

# Robustness against adversarial attacks

- Similarities of EBM training to Adversarial training:
  - We use a gradient-based optimization to generate examples which activate a high-level network activation. Then optimize the network to minimize the generated example's effect on the activation.
  - Given these connections one may wonder if JEM would be more robust to adversarial examples than a standard model. This behavior has been demonstrated in prior work on EBMs.
- Adversarial Attacks [13]
  - Adversarial threats are modeled as perturbed inputs:  $\tilde{x} = x + \delta$  where  $\|\tilde{x} - x\|_p < \varepsilon$
  - Closeness with regards to  $L_p$  distance does not imply that adversarial examples reside within areas of high density according to the model distribution.



# Robustness against adversarial attacks

- Perturbation-based adversarial examples
  - We propose to run a few iterations of the sampling procedure seeded at the given input. This should be able to transform low-probability inputs to a nearby point of high probability, “undoing” any adversarial attack and enabling the model to classify robustly.
  - We test JEM with 0, 1, 10 steps of sampling seeded at the input.
  - Our model is considerably more robust than a baseline with standard classifier training. JEM with 0 steps refinement is noticeably more robust than the baseline.
  - Adv and RandAdvSmooth [14] is trained to be robust to the specific norm. However, we attack the same JEM model with both norms and observe competitive robustness in both cases.

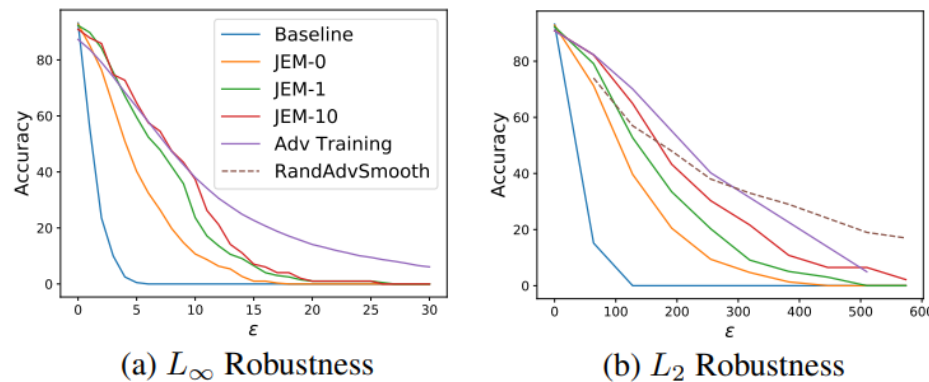


Figure 5: Adversarial Robustness Results with PGD attacks. JEM adds considerable robustness.

# Robustness against adversarial attacks

- Distal adversarial examples [15]
  - Non-robust models have tendency to classify non-sensical inputs with high confidence.
  - Starting from noise we generate images to maximize  $p_{\theta}(y = \text{car}|\mathbf{x})$ .
  - The baseline confidently classifies unstructured noise images.
  - The L2-adversarially trained ResNet with  $\epsilon = 0.5$  confidently classifies somewhat structured, but unrealistic images.
  - JEM does not confidently classify nonsensical images, instead, car attributes and natural image properties visibly emerge.

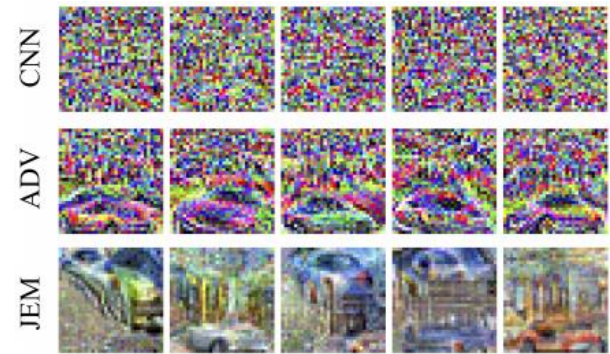


Figure 6: **Distal Adversaries.** Confidently classified images generated from noise, such that:  $p(y = \text{"car"}|\mathbf{x}) > .9$ .

# Limitations

- EBM's can be very challenging to work with:
  - Difficulty in evaluation: Since normalized likelihoods cannot be computed, it is hard to even verify that learning is taking place at all. In some domains such as images, samples could be drawn to assess learning, however this is far from a generalizable strategy.
  - Difficulty in training: The gradient estimators we use to train JEM are quite unstable and are prone to diverging if the sampling and optimization parameters are not tuned correctly. The models demonstrated in this work regularly diverged throughout training, requiring them to be restarted with lower learning rates or with increased regularization.
- While this may seem prohibitive, we believe that the results presented in this work are sufficient to motivate the community to find solutions to these issues as any improvement in the training of energy based models will further improve the results we have presented in this work.

# Conclusion

- We have presented JEM, a novel reinterpretation of standard classifier architectures which retains the strong performance of SOTA discriminative models while adding the benefits of generative modeling approaches.
- Our work is based upon recent works on scaling techniques for training EBMs to high-dimensional data.
- We have demonstrated the utility of incorporating this type of training into discriminative models.



# References

- Will Grathwohl, Kuan-Chieh Wang, Joern-Henrik Jacobsen, David Duvenaud, Mohammad Norouzi, Kevin Swersky. Your classifier is secretly an energy based model and you should treat it like one. In *International Conference on Learning Representations (ICLR)*, 2020