

Conjugate Energy-Based Models

Hao Wu*

Northeastern University

WU.HAO10@NORTHEASTERN.EDU

Babak Esmaeili*

Northeastern University

ESMAEILI.B@NORTHEASTERN.EDU

Michael Wick

Oracle Labs

MICHAEL.WICK@ORACLE.COM

Jean-Baptiste Tristan

Boston College

TRISTANJ@BC.EDU

Jan-Willem van de Meent

Northeastern University

J.VANDEMEENT@NORTHEASTERN.EDU

Abstract

We propose conjugate energy-based models (EBMs), a class of deep latent-variable models with a tractable posterior. Conjugate EBMs have similar use cases as variational autoencoders, in the sense that they learn an unsupervised mapping between data and latent variables. However these models omit a generator, which allows them to learn more flexible notions of similarity between data points. Our experiments demonstrate that conjugate EBMs achieve competitive results in terms of image modelling, predictive power of latent space, and out-of-distribution detection on a variety of datasets.

1. Introduction

Deep generative models approximate a data distribution by combining a prior over latent variables with a neural generator, which maps latent variables to points on a data manifold. While it is common to evaluate these models in terms of their ability to generate realistic examples, an arguably more important use case of these models is that they learn representations in absence of supervision. To be useful in downstream tasks, these representations should encode some set of “semantically meaningful” features rather than “nuisance variables” that are unlikely to have predictive power.

Guiding a model towards a semantically meaningful representation requires some form of inductive bias. A large body of work on variational autoencoders (VAEs, [Kingma and Welling \(2013\)](#); [Rezende et al. \(2014\)](#)) has explored the use of priors as inductive biases. Relatively mild biases in the form of conditional independence are common in the literature on disentangled representations ([Higgins et al., 2016](#); [Kim and Mnih, 2018](#); [Chen et al., 2018](#); [Esmaeili et al., 2019](#)). More generally, recent work has employed priors that reflect structure of the underlying data to represent objects in an image ([Eslami et al., 2016](#); [Lin et al., 2020b](#); [Engelcke et al., 2019](#); [Crawford and Pineau, 2019b](#)), or moving objects in video ([Crawford and Pineau, 2019a](#); [Kosiorrek et al., 2018](#); [Wu et al., 2020](#); [Lin et al., 2020a](#)).

Despite steady progress, work on disentangled representations and structured VAEs still predominantly considers synthetic data sets. To train a VAE we minimize a reconstruction loss, which treats all pixels in an image equally. For complex natural scenes, learning a

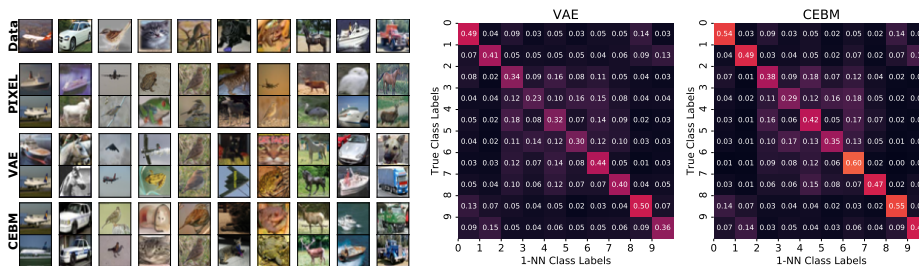


Figure 1: (*Left*) Samples from CIFAR-10 along with the top 2-nearest-neighbors in pixel space, the latent space of a VAE, and the latent space of a CEBM. (*Right*) Confusion matrices of 1-nearest-neighbor classification on CIFAR-10 based on L2 distance in the latent space. On average, CEBM representations more closely align with class labels compared to VAE.

model that can produce pixel-perfect reconstructions poses fundamental challenges, given the combinatorial explosion of possible inputs. This is not only a problem from the perspective of generation, but also from the perspective of the learned representation; a VAE must encode all factors of variation that give rise to large deviations in pixel space, regardless of whether these factors are semantically meaningful (e.g. presence and locations of objects) or not (e.g. shadows of objects in the background of the image).

In this paper, we consider energy-based models (EBMs) with latent variables as an alternative VAEs for learning representations in an unsupervised manner. The general idea of using EBMs for this purpose is by no means new; it has a long history in the context of restricted Boltzmann machines (RBMs) and related models (Smolensky, 1986; Welling et al., 2004; Hinton et al., 2006). Our motivation with the present work is to design a class of energy-based models that retain the desirable features of VAEs while addressing what we see as one of their main weaknesses: We would like incorporate inductive biases, but model the data at an intermediate level of representation that does not necessarily encode all features of an image at the pixel level.

Concretely, we propose Conjugate EBMs (CEBMs), a class of models in which the energy function defines a neural exponential family. While the normalizer of this family is intractable, we can compute its posterior in closed form when we pair the likelihood with an appropriate conjugate bias term in the energy function. As a result, the neural sufficient statistics in a CEBM fully determine both the marginal likelihood and the encoder, hereby side-stepping the need for a generator.

In our experiments, we evaluate the representations learned by CEBM using class labels as a proxy for the primary factors of variation in a dataset. We show that CEBMs learn a notion of similarity that aligns more closely with class labels in terms of the nearest neighbors in latent space (Figure 1). Moreover, we show that the representations learned by CEBMs (in an unsupervised manner) can achieve a competitive performance in classification and out-of-distribution detection tasks.

2. Background

2.1. Energy-Based Models

An EBM (LeCun et al., 2006) defines a density $p_\theta(x) = \exp\{-E_\theta(x)\}/Z_\theta$ in terms of an energy function $E_\theta(x) : \mathbb{R}^D \rightarrow \mathbb{R}$ with parameters θ , which maps inputs $x \in \mathbb{R}^D$ to a scalar value. The density $p_\theta(x)$ can only be evaluated up to an unknown constant of proportionality since the normalizing constant $Z_\theta = \int dx \exp\{-E_\theta(x)\}$ is typically intractable. To train an EBM, we approximate the gradient of the expected log-likelihood

$$\nabla_\theta \mathcal{L}_\theta = \mathbb{E}_{p_{\text{data}}(x)}[\nabla_\theta \log p_\theta(x)] = \mathbb{E}_{p_{\text{data}}(x)}[-\nabla_\theta E_\theta(x)] - \nabla_\theta \log Z_\theta, \quad (1)$$

$$= \mathbb{E}_{p_{\text{data}}(x)}[\nabla_\theta E_\theta(x)] + \mathbb{E}_{p_\theta(x')}[\nabla_\theta E_\theta(x')]. \quad (2)$$

Estimating the gradient of the log normalizer $\nabla_\theta \log Z_\theta$ requires samples from the model $x' \sim p_\theta(x')$, which in turn requires approximate inference since this density is intractable. A commonly used method is Stochastic Gradient Langevin Dynamics (SGLD, Welling and Teh (2011)), which initializes a sample $x'_0 \sim p_0(x')$ and then performs a sequence of gradient updates with additional injected noise ϵ

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

Recent work has shown that EBMs with convolutional energy functions can accurately model distributions over images, in the sense that SGLD produces realistic samples (Nijkamp et al., 2019a,b; Du and Mordatch, 2019; Xie et al., 2016).

2.2. Conjugate Exponential Families

An exponential family is a set of distributions whose probability density function or probability mass function can be expressed in the form

$$p(x | \eta) = h(x) \exp\{t(x)^\top \eta - A(\eta)\},$$

where $h(\cdot)$ is base measure, η is a vector of natural parameters, $t(\cdot)$ is a vector of sufficient statistics, and $A(\cdot)$ is a log normalizer. If a likelihood belongs to an exponential family, then there exists a conjugate prior with the form $p(\eta | \lambda) = h_0(\eta) \exp\{t_0(\eta)^\top \lambda - A_0(\lambda)\}$. A prior is conjugate to a likelihood when its vector of sufficient statistics comprises the natural parameters and the log-normalizer of the likelihood $t_0(\eta) = [\eta, -A(\eta)]$, $\lambda = [\lambda_1, \lambda_2]$. The convenient property of conjugate exponential families is that both the marginal likelihood $p(x | \lambda)$ and the posterior $p(\eta | x, \lambda)$ are tractable. The reason is that the joint probability has the form

$$p(x, \eta | \lambda) = p(x | \eta) p(\eta | \lambda) = \exp\{\eta^\top (\lambda_1 + t(x)) - A(\eta)(\lambda_2 + 1) - A_0(\lambda)\}.$$

If we substitute $\tilde{\lambda}_1 = \lambda_1 + t(x)$ and $\tilde{\lambda}_2 = \lambda_2 + 1$, we can equivalently factorize this joint as

$$p(x, \eta | \lambda) = p(\eta | x, \lambda) p(x | \lambda) = p(\eta | \tilde{\lambda}) \exp\{A_0(\tilde{\lambda}) - A_0(\lambda)\}. \quad (3)$$

This shows that the posterior is in the same exponential family as the prior, and that we can express the marginal likelihood using the log normalizer $A_0(\cdot)$

$$p(\eta | x, \lambda) = p(\eta | \tilde{\lambda}), \quad p(x | \lambda) = \exp\{A_0(\tilde{\lambda}) - A_0(\lambda)\}. \quad (4)$$

3. Conjugate Energy-Based Models

In this paper we are interested in learning a probabilistic model that defines a joint distribution $p_\theta(x, z)$ over high-dimensional data $x \in \mathbb{R}^D$ and a lower-dimensional set of latent variables $z \in \mathbb{R}^K$. Because most work on deep generative models has focused on images, we will restrict ourselves to this data modality. We propose to consider models of the form

$$p_\theta(x, z) = \frac{1}{Z_\theta} \exp \{ -E_\theta(x, z) \}, \quad E_\theta(x, z) = -t_\theta(x)^\top \eta(z) - b(z), \quad (5)$$

In the energy function, $\eta : \mathbb{R}^K \rightarrow \mathbb{R}^H$ decodes latent variables to a vector of natural parameters in an intermediate space of dimension H . The function $t_\theta : \mathbb{R}^D \rightarrow \mathbb{R}^H$ plays the role of an encoder, which maps data to sufficient statistics in the same space as $\eta(z)$. The function $b : \mathbb{R}^K \rightarrow \mathbb{R}$ serves as an inductive bias that plays a role analogous to the prior. We will consider a bias $b(z) = \log p(z | \lambda)$ in form of a tractable exponential family

$$b(z) = \eta(z)^\top \lambda - A(\lambda). \quad (6)$$

We can then express the energy function as

$$E_\theta(x, z) = -t_\theta(x)^\top \eta(z) - b(z) = -(\lambda + t_\theta(x))^\top \eta(z) + A(\lambda). \quad (7)$$

This form of the energy function has a very convenient property: It corresponds to a model $p_\theta(x, z)$ in which the posterior $p_\theta(z | x)$ is tractable. To see this, we can make a substitution $\tilde{\lambda} = \lambda + t_\theta(x)$ analogous to the one in Equation 3, which allows us to express the energy as

$$E_\theta(x, z) = -(\eta(z)^\top \tilde{\lambda} - A(\tilde{\lambda})) - (A(\tilde{\lambda}) - A(\lambda)),$$

We now see that we can factorize the corresponding density $p_\theta(x, z)$ as

$$p_\theta(x, z | \lambda) = p_\theta(x | \lambda) p_\theta(z | x, \lambda), \quad (8)$$

which yields a posterior and a marginal that are analogous to the distributions in Equation 4

$$p_\theta(z | x, \lambda) = p(z | \lambda + t_\theta(x)), \quad p_\theta(x | \lambda) = \frac{1}{Z_\theta} \exp \{ A(\lambda + t_\theta(x)) - A(\lambda) \}. \quad (9)$$

This posterior is conjugate, in the sense that it is in the same tractable family as the bias. In addition to having a tractable posterior, CEBMs have the convenient property that the marginal likelihood $p_\theta(x | \lambda)$ itself can be expressed as an energy-based model that is defined in terms of the log normalizer $A(\cdot)$ and the encoder network $t_\theta(x)$. This means that we can train CEBMs using SGLD in the same way as other EBMs.

CEBMs differ from VAEs in that they lack a generator network. Instead, the density is fully specified by the encoder network $t_\theta(x)$, which defines a notion of agreement $(\lambda + t_\theta(x))^\top \eta(z)$ between data and latent variables in an intermediate feature space. In VAEs, by contrast we would assume that the sufficient statistics $t(x)$ are known, and learn a generator $\eta_\theta(z)$ to compute a notion of agreement in data space (see Appendix B). Note that when z is a categorical variable, the sufficient statistics $t_\theta(x)$ define a classifier, and CEBMs are equivalent to the models that are considered in recent work on EBMs for classification (Grathwohl et al., 2019; Liu and Abbeel, 2020).



Figure 2: CEBM Samples of MNIST, F-MNIST, SVHN and CIFAR-10.

Models	MNIST				Fashion-MNIST				CIFAR-10				SVHN			
	1	10	100	full	1	10	100	full	1	10	100	full	1	10	100	full
VAE	42	85	92	95	41	63	72	81	16	22	31	38	13	13	16	36
GMM-VAE	53	86	93	97	49	68	79	84	19	23	33	39	13	14	23	56
IGEBM	63	89	95	97	50	70	79	83	16	26	33	42	10	16	35	49
CEBM	67	89	95	97	52	70	77	83	19	30	42	52	12	25	48	70
CEBMM	67	91	97	98	52	71	80	85	16	28	42	51	10	17	39	60

Table 1: Classification accuracy. We pre-train 5 unsupervised models (rows) on MNIST, Fashion-MNIST, CIFAR10, SVHN. Then we train logistic classifiers using 1, 10, 100 examples per class (i.e. shots) and the full training dataset. We report the testing classification accuracy, where CEBM outperforms.

4. Inductive Biases

CEBMs have a property that is somewhat counter-intuitive. While the posterior $p_\theta(z | x, \lambda)$ is tractable, the prior $p_\theta(z)$ is in general not tractable. In particular, although the bias $b_\theta(z)$ is the logarithm of a tractable exponential family, it is not the case that $p_\theta(z) \neq \exp\{b_\theta(z)\}$. Rather the prior $p_\theta(z)$ has the form,

$$p_\theta(z) = \frac{\exp\{b_\theta(z)\}}{Z_\theta} \int dx \exp\{t_\theta(x)^\top \eta(z)\}.$$

In principle the bias a CEBM can take the form of any exponential family distribution. Since products of exponential families are also in the exponential family, this covers a broad range of possible biases. In this paper, we will constrain ourselves to a Spherical Gaussian and a Mixture of Gaussians. We provide derivations for both cases in Appendix A.

5. Experiments

Our experiments evaluate to what extent CEBMs can learn representations that encode meaningful factors of variation, whilst discarding details about the input that we would consider noise. We train using both inductive biases and will refer to them as CEBM (Spherical Gaussian) and CEBMM (Mixture of Gaussian) for the rest of Section. See Appendix C and D for architecture and training details.

5.1. Samples and Latent Space

We begin with a qualitative evaluation by visualizing samples. Figure 2 shows samples from CEBMs with uniform noise initialization and 500 SGLD steps.

To assess to what extent the representation in CEBMs aligns with classes in each dataset, we look at the agreement between the label for each data point and the label of its nearest neighbor in the latent space. Figure 1 shows that the distance in pixel space is a poor measure of similarity in this dataset, whereas proximity in the latent space is more likely to agree with class labels in both VAEs and CEBMs. Quantitatively, in Appendix F, we observe a stronger alignment between classes and the latent representation in CEBMs, which is reflected in higher numbers on the diagonal of the matrix.

5.2. Classification

To evaluate performance in settings where few labels are available, we train a logistic classifier using 1, 10, 100 examples per class, as well as the full training dataset. We compare CEBMs against the IGBM (Du and Mordatch, 2019), a standard VAE with the spherical Gaussian prior, and the GMM-VAE (Tomczak and Welling, 2018) where the prior is a mixture of Gaussians. As discussed in Section 2, IGBM does not have an explicit representation. In order to compare against IGBM, we remove the last layer (which outputs the energy) and use the resulting intermediate representation as the latent code.

We report the classification accuracy on the test set in Table 1. We can see that that CEBMs overall achieve a higher accuracy compared to VAEs in particular for CIFAR-10 and SVHN where the pixel distance is not good measure for similarity. Moreover, we observe that CEBMs outperform IGBM which suggest that the inductive biases in CEBMs can lead to increased performance in downstream tasks.

5.3. Out-of-Distribution Detection

EBMs have formed the basis for encouraging results in out-of-distribution (OOD) detection (Du and Mordatch, 2019; Grathwohl et al., 2019). In Appendix G, we report the area under the receiver-operator curve (AUROC) using two score functions: $\log p_{\theta}(x)$ and a gradient-based function proposed by Grathwohl et al. (2019). CEBMs results for OOD detection in most cases improve upon VAE and IGBM baselines.

6. Conclusion

We introduced CEBMs; a new family of energy-based models that define a joint energy function over both the data and latent variables. The joint distribution factorizes into a tractable posterior and a marginal likelihood, imposing an inductive bias on the latent space. This factorization allows us to directly optimize the marginal likelihood of the data, while at the same time imposing an inductive bias on the latent space. Experimental results for this class of models are encouraging; we observe a closer agreement between unsupervised representations and class labels, which translates into improvements in downstream classification tasks.

Acknowledgments

We would like to thank our reviewers for their thoughtful comments. This work was supported by the Intel Corporation, the 3M Corporation, NSF award 1835309, startup funds from Northeastern University, the Air Force Research Laboratory (AFRL), and DARPA.

References

- Ricky TQ Chen, Xuechen Li, Roger B Grosse, and David K Duvenaud. Isolating sources of disentanglement in variational autoencoders. In *Advances in neural information processing systems*, pages 2610–2620, 2018.
- Eric Crawford and Joelle Pineau. Exploiting spatial invariance for scalable unsupervised object tracking. *arXiv preprint arXiv:1911.09033*, 2019a.
- Eric Crawford and Joelle Pineau. Spatially invariant unsupervised object detection with convolutional neural networks. In *Proceedings of the AAAI Conference on Artificial Intelligence*, volume 33, pages 3412–3420, 2019b.
- Yilun Du and Igor Mordatch. Implicit generation and generalization in energy-based models. *arXiv preprint arXiv:1903.08689*, 2019.
- Martin Engelcke, Adam R Kosior, Oivi Parker Jones, and Ingmar Posner. Genesis: Generative scene inference and sampling with object-centric latent representations. *arXiv preprint arXiv:1907.13052*, 2019.
- S. M. Ali Eslami, Nicolas Heess, Theophane Weber, Yuval Tassa, David Szepesvari, Koray Kavukcuoglu, and Geoffrey E. Hinton. Attend, infer, repeat: Fast scene understanding with generative models. In *Proceedings of the 30th International Conference on Neural Information Processing Systems, NIPS’16*, pages 3233–3241, Red Hook, NY, USA, December 2016. Curran Associates Inc. ISBN 978-1-5108-3881-9.
- Babak Esmaeili, Hao Wu, Sarthak Jain, Alican Bozkurt, Narayanaswamy Siddharth, Brooks Paige, Dana H Brooks, Jennifer Dy, and Jan-Willem Meent. Structured disentangled representations. In *The 22nd International Conference on Artificial Intelligence and Statistics*, pages 2525–2534. PMLR, 2019.
- Will Grathwohl, Kuan-Chieh Wang, Jörn-Henrik Jacobsen, David Duvenaud, Mohammad Norouzi, and Kevin Swersky. Your classifier is secretly an energy based model and you should treat it like one. *arXiv preprint arXiv:1912.03263*, 2019.
- Irina Higgins, Loic Matthey, Arka Pal, Christopher Burgess, Xavier Glorot, Matthew Botvinick, Shakir Mohamed, and Alexander Lerchner. beta-vae: Learning basic visual concepts with a constrained variational framework. 2016.
- Geoffrey E. Hinton, Simon Osindero, and Yee-Whye Teh. A Fast Learning Algorithm for Deep Belief Nets. *Neural Computation*, 18(7):1527–1554, May 2006. ISSN 0899-7667. doi: 10.1162/neco.2006.18.7.1527.
- Hyunjik Kim and Andriy Mnih. Disentangling by factorising. In *International Conference on Machine Learning*, pages 2649–2658, 2018.
- Diederik P. Kingma and Max Welling. Auto-encoding variational bayes. *International Conference on Learning Representations*, 2013.

- Adam Kosior, Hyunjik Kim, Yee Whye Teh, and Ingmar Posner. Sequential attend, infer, repeat: Generative modelling of moving objects. In *Advances in Neural Information Processing Systems*, pages 8606–8616, 2018.
- Yann LeCun, Sumit Chopra, Raia Hadsell, M Ranzato, and F Huang. A tutorial on energy-based learning. *Predicting structured data*, 1(0), 2006.
- Zhixuan Lin, Yi-Fu Wu, Skand Peri, Bofeng Fu, Jindong Jiang, and Sungjin Ahn. Improving generative imagination in object-centric world models. *arXiv preprint arXiv:2010.02054*, 2020a.
- Zhixuan Lin, Yi-Fu Wu, Skand Vishwanath Peri, Weihao Sun, Gautam Singh, Fei Deng, Jindong Jiang, and Sungjin Ahn. SPACE: Unsupervised Object-Oriented Scene Representation via Spatial Attention and Decomposition. *arXiv:2001.02407 [cs, eess, stat]*, March 2020b.
- Hao Liu and Pieter Abbeel. Hybrid discriminative-generative training via contrastive learning. *arXiv preprint arXiv:2007.09070*, 2020.
- Erik Nijkamp, Mitch Hill, Tian Han, Song-Chun Zhu, and Ying Nian Wu. On the anatomy of mcmc-based maximum likelihood learning of energy-based models. *arXiv*, pages arXiv-1903, 2019a.
- Erik Nijkamp, Mitch Hill, Song-Chun Zhu, and Ying Nian Wu. Learning non-convergent non-persistent short-run mcmc toward energy-based model. In *Advances in Neural Information Processing Systems*, pages 5232–5242, 2019b.
- Danilo Jimenez Rezende, Shakir Mohamed, and Daan Wierstra. Stochastic Backpropagation and Approximate Inference in Deep Generative Models. In Eric P. Xing and Tony Jebara, editors, *Proceedings of the 31st International Conference on Machine Learning*, volume 32 of *Proceedings of Machine Learning Research*, pages 1278–1286, Beijing, China, June 2014. PMLR.
- Paul Smolensky. Information processing in dynamical systems: Foundations of harmony theory. Technical report, Colorado Univ at Boulder Dept of Computer Science, 1986.
- Jakub Tomczak and Max Welling. Vae with a vampprior. In *International Conference on Artificial Intelligence and Statistics*, pages 1214–1223, 2018.
- Max Welling and Yee W Teh. Bayesian learning via stochastic gradient langevin dynamics. In *Proceedings of the 28th international conference on machine learning (ICML-11)*, pages 681–688, 2011.
- Max Welling, Michal Rosen-zvi, and Geoffrey E Hinton. Exponential family harmoniums with an application to information retrieval. In L. Saul, Y. Weiss, and L. Bottou, editors, *Advances in Neural Information Processing Systems*, volume 17, pages 1481–1488. MIT Press, 2004.

Ronald J. Williams. Simple statistical gradient-following algorithms for connectionist reinforcement learning. *Machine Learning*, 8(3-4):229–256, 1992. ISSN 0885-6125. doi: 10.1007/BF00992696.

Hao Wu, Heiko Zimmermann, Eli Sennesh, Tuan Anh Le, and Jan-Willem van de Meent. Amortized population gibbs samplers with neural sufficient statistics. In *Proceedings of the International Conference on Machine Learning*, pages 10205–10215. 2020.

Jianwen Xie, Yang Lu, Song-Chun Zhu, and Yingnian Wu. A theory of generative convnet. In *International Conference on Machine Learning*, pages 2635–2644, 2016.

Appendix A. Derivation for Two Cases of Inductive Biases

1. Spherical Gaussian. As a bias that is analogous to the standard prior in VAEs, we consider a spherical Gaussian with fixed hyperparameters $(\mu, \sigma) = (0, 1)$ for each dimension of $z \in \mathbb{R}^K$,

$$b_\theta(z) = \sum_k (\eta(z_k)^\top \lambda - A(\lambda)),$$

Each term has sufficient statistics $\eta(z_k) = (z_k, z_k^2)$ and natural parameters

$$\lambda = \left(\frac{\mu}{\sigma^2}, -\frac{1}{2\sigma^2} \right) = \left(0, -\frac{1}{2} \right). \quad (10)$$

The marginal likelihood of the CEBM is then

$$p_\theta(x | \lambda) = \frac{1}{Z_\theta} \exp \left\{ \sum_k (A(\tilde{\lambda}_k) - A(\lambda)) \right\}, \quad (11)$$

where $\tilde{\lambda}_k = \lambda + t_{\theta,k}(x)$ and the log normalizer is

$$A(\lambda_k) = -\frac{\lambda_{k,1}^2}{4\lambda_{k,2}} - \frac{1}{2} \log(-\lambda_{k,2}).$$

2. Mixture of Gaussians. In our experiments we will consider datasets that are normally used for classification. These datasets, by design, exhibit multimodal structure that we would like to see reflected in the learned representation. As an inductive bias that is amenable to uncovering this structure, we will consider a bias in the form of a mixture of L Gaussians,

$$b_\theta(y, z) = \sum_{k,l} I[y = l] (\eta(z_k)^\top \lambda_{l,k} - A(\lambda_{l,k})).$$

Here $z \in \mathbb{R}^K$ is a vector of features and $y \in \{1, \dots, L\}$ is a categorical assignment variable. The bias for each component l is a spherical Gaussian with hyperparameters $\lambda_{l,k}$ for each dimension k . Again using the notation $\tilde{\lambda}_{l,k} = \lambda_{l,k} + t_{\theta,l,k}(x)$ to refer to the posterior parameters, then we obtain an energy

$$E_\theta(x, y, z) = -\sum_{k,l} I[y = l] (\tilde{\lambda}_{l,k}^\top \eta(z_k) - A(\lambda_{l,k})).$$

We can then define a joint probability over data x and the assignment y in terms the log normalizer $A(\cdot)$,

$$p_\theta(x, y | \lambda) = \frac{1}{Z_\theta(\lambda)} \exp \left\{ \sum_{k,l} I[y = l] (A(\tilde{\lambda}_{l,k}) - A(\lambda_{l,k})) \right\},$$

which then allows us to compute the marginal

$$p_\theta(x | \lambda) = \sum_y p_\theta(x, y | \lambda). \quad (12)$$

We optimize this marginal with respect hyperparameters λ as well as the weights θ .

Appendix B. Relationship between CEBMs and VAEs

B.1. Variational Autoencoders

Variational autoencoders are a widely used class of deep generative models [Kingma and Welling \(2013\)](#); [Rezende et al. \(2014\)](#). A VAE defines a joint distribution $p_\theta(x, z)$ over data x and latent variables z ; it combines an unstructured prior (e.g. a spherical Gaussian) with a likelihood that is parameterized by an expressive neural network, often referred to as a decoder. An inference model, also known as an encoder, approximates the posterior $p_\theta(z | x)$ by mapping each data point x onto latent variables z . These models are trained by maximizing the stochastic evidence lower bound (ELBO) defined as

$$\mathcal{L}(\phi, \theta) = \mathbb{E}_{p_{\text{data}}(x) q_\phi(z|x)} \left[\log \frac{p_\theta(x, z)}{q_\phi(z | x)} \right] \quad (13)$$

When the $q_\phi(z|x)$ is reparameterizable, we can compute Monte Carlo estimates of the gradient of this objective using pathwise derivatives. Non-reparameterizable cases, such as models with discrete variables, require likelihood-ratio estimators [Williams \(1992\)](#).

Despite their successes, VAEs have limitations. By maximizing the ELBO, we favor encoder-decoder pair that perfectly reconstruct their input, so there is nothing preventing the VAE from mapping similar inputs to similar encoding, even when they might be semantically different. Likewise, examples that might be very dissimilar in pixel space because of noise or benign transformations might end up with very different latent representations.

B.2. VAEs as Energy-based Models

We can interpret the generative model in a VAE as a model with an energy function

$$E_\theta(x, z) = -t(x)^\top \eta_\theta(z) - b_\theta(z). \quad (14)$$

In this setting, $\eta_\theta(z)$ is the generator network that maps low-dimensional latent variables to a high-dimensional vector of natural parameters. The function $t(x)$ is a known mapping of data to the sufficient statistics of a Gaussian or Bernoulli likelihood.

The bias $b_\theta(z)$ contains the terms in the log density $\log p_\theta(x, z)$ that only depend on z . In a VAE there are two such terms. The first is the log prior $\log p_\theta(z)$. The second is the log normalizer of $A(\eta_\theta(z))$ for the likelihood $\log p_\theta(x | z)$, which can be computed in closed form because Gaussian and Bernoulli distributions are in the exponential family. Combining these terms yields an energy function for the bias,

$$b_\theta(z) = \log p_\theta(z) - A(\eta_\theta(z)). \quad (15)$$

In other words, in a VAE we use known sufficient statistics $t(x)$, and train a generator to learn the natural parameters $\eta_\theta(z)$. In a CEBM, by contrast, we assume known natural parameters $\eta(z)$ and train an encoder to learn the sufficient statistics $t_\theta(x)$.

Appendix C. Model Architectures

CEBMs employ an encoder network $t_\theta(x)$ in the form of 4-layer CNN (which is proposed by Nijkamp et al. (2019a)), followed by an MLP output layer. For IGEBMs, we add one extra MLP as its final layer which outputs a scalar value. VAEs use the same encoder network as the CEBMs, and use a decoder network in form of an MLP followed by 4-layer CNN.

Table Table C, Table Table C, and Table Table C show the architectures used for CEBM, VAE, and IGEBM, respectively.

(a) MNIST and Fashion-MNIST.	(b) CIFAR10 and SVHN.
Encoder	Encoder
Input $28 \times 28 \times 1$ images	Input $32 \times 32 \times 3$ images
3×3 conv. 64 LeakyReLU. stride 1. padding 1	3×3 conv. 64 LeakyReLU. stride 1. padding 1
4×4 conv. 64 LeakyReLU. stride 2. padding 1	4×4 conv. 128 LeakyReLU. stride 2. padding 1
4×4 conv. 32 LeakyReLU. stride 2. padding 1	4×4 conv. 256 LeakyReLU. stride 2. padding 1
4×4 conv. 32 LeakyReLU. stride 2. padding 1	4×4 conv. 512 LeakyReLU. stride 2. padding 1
FC. 128 LeakyReLU	FC. 128 LeakyReLU
FC. 2×128	FC. 2×128

Table 2: Architecture of CEBM.

(a) MNIST and Fashion-MNIST.	
Encoder	Decoder
Input $28 \times 28 \times 1$ images	Input $z \in \mathbb{R}^{128}$ latent variables
3×3 conv. 64 LeakyReLU. stride 1. padding 1	FC. 128 ReLU
4×4 conv. 64 LeakyReLU. stride 2. padding 1	FC. $3 \times 3 \times 32$ ReLU
4×4 conv. 32 LeakyReLU. stride 2. padding 1	4×4 upconv. 32 LeakyReLU. stride 2. padding 1
4×4 conv. 32 LeakyReLU. stride 2. padding 1	4×4 upconv. 64 LeakyReLU. stride 2. padding 1
FC. 128 ReLU	4×4 upconv. 64 LeakyReLU. stride 2. padding 0
FC. 2×128	3×3 upconv. 1 stride 1. padding 0
(b) CIFAR10 and SVHN.	
Encoder	Decoder
Input $32 \times 32 \times 3$ images	Input $z \in \mathbb{R}^{128}$ latent variables
3×3 conv. 64 LeakyReLU. stride 1. padding 1	FC. 128 ReLU
4×4 conv. 128 LeakyReLU. stride 2. padding 1	FC. $4 \times 4 \times 512$ ReLU
4×4 conv. 256 LeakyReLU. stride 2. padding 1	4×4 upconv. 32 LeakyReLU. stride 2. padding 1
4×4 conv. 512 LeakyReLU. stride 2. padding 1	4×4 upconv. 64 LeakyReLU. stride 2. padding 1
FC. 128 ReLU	3×3 upconv. 64 LeakyReLU. stride 2. padding 1
FC. 2×128	3×3 upconv. 1 stride 1. padding 1

Table 3: Architecture of VAE.

(a) MNIST and Fashion-MNIST.	(b) CIFAR10 and SVHN.
Encoder	Encoder
Input $28 \times 28 \times 1$ images	Input $32 \times 32 \times 3$ images
3×3 conv. 64 LeakyReLU. stride 1. padding 1	3×3 conv. 64 LeakyReLU. stride 1. padding 1
4×4 conv. 64 LeakyReLU. stride 2. padding 1	4×4 conv. 128 LeakyReLU. stride 2. padding 1
4×4 conv. 32 LeakyReLU. stride 2. padding 1	4×4 conv. 256 LeakyReLU. stride 2. padding 1
4×4 conv. 32 LeakyReLU. stride 2. padding 1	4×4 conv. 512 LeakyReLU. stride 2. padding 1
FC. 128 LeakyReLU	FC. 128 LeakyReLU
FC. 128 LeakyReLU. FC. 1	FC. 128 LeakyReLU. FC. 1

Table 4: Architecture of IGEBM

Appendix D. Training Details of Persistent Contrastive Divergence

Optimization. In CEBMs and VAEs, we choose the dimension of latent variables to be 128. For CEBMs, We found that the optimization becomes difficult with smaller dimensions. We L2 regularize energy magnitudes (proposed by [Du and Mordatch \(2019\)](#)), where the coefficient of the L2 regularization term is 0.1. We empirically found that the training would become unstable without this regularization. We train our models using 60 SGLD steps where we initialize samples from the replay buffer with 0.95 probability, and initialize from uniform noise with 0.05 probability. We train all the models with 90k gradient steps, batch size 128, Adam optimizer with learning rate 1e-4. When doing PCD, we used a replay buffer of size 5000. We set the α in the SGLD steps to be 0.075. Similar to [Du and Mordatch \(2019\)](#), we found it useful to add some noise to the image before encoding. In our experiments, we used Gaussian noise with $\sigma^2 = 0.03$. For the mixture models (CEBMM and GMM-VAE), we used 50 mixtures.

Training Stability. As observed in previous work [Du and Mordatch \(2019\)](#); [Grathwohl et al. \(2019\)](#), training EBMs can be a challenging task that often requires a thorough hyperparameter search. We found that the choices of activation function, learning rate, number of SGLD steps, and regularization will all affect training stability. Models regularly diverge during training, and it is difficult to perform diagnostics given that $\log p_\theta(x)$ cannot be computed. As suggested by [Nijkamp et al. \(2019a\)](#), we found checking the difference in energy between data and model samples to be helpful for verifying stable training. We also note that in general, we observed a trade-off between sample quality and the predictive power of latent variables in our experiments. We leave investigation of the source of this trade-off to future work, but we suspect that this is because SGLD is having more difficulty to converge when the latent space is more disjoint.

Appendix E. Persistent Contrastive Divergence

Algorithm 1: Persistent Contrastive Divergence

Input: $p_{\text{data}}(\cdot)$, θ , α , T
 $\mathcal{B} \leftarrow \{x_b \sim \mathcal{U} \text{ for } b = 1 \dots \text{buffer-size}\};$
while *not converged* **do**
 $x \sim p_{\text{data}}(x);$
 $x' \sim \mathcal{B}$ with 95% probability and \mathcal{U} otherwise;
 for $t = 1 \dots T$ **do**
 $\epsilon \sim \mathcal{N}(0, \alpha);$
 $x' \leftarrow x' - \frac{\alpha}{2} \nabla_x E_{\theta}(x') + \epsilon;$
 end
 $\Delta_{\theta} \leftarrow \nabla_{\theta} E_{\theta}(x) - \nabla_{\theta} E_{\theta}(x');$
 $\theta \leftarrow \text{Adam}(\theta, \Delta_{\theta});$
 $\mathcal{B}[x'] \leftarrow x';$
end
Output: θ

Appendix F. Confusion Matrices on 1-NN Classification

We perform 1-nearest-neighbor classification task for MNIST, Fashion-MNIST, SVHN, CI-FAR10. We compute the L2 distance in the latent space of VAE, IGEBM and CEBM, and also in pixel space. We visualize the confusion matrices

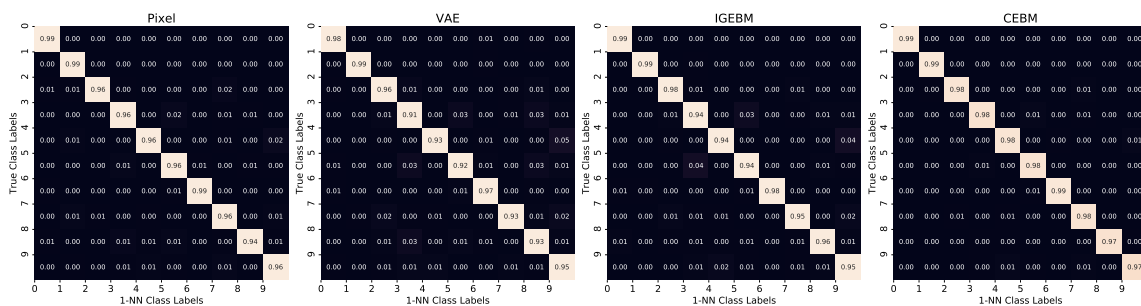


Figure 3: MNIST

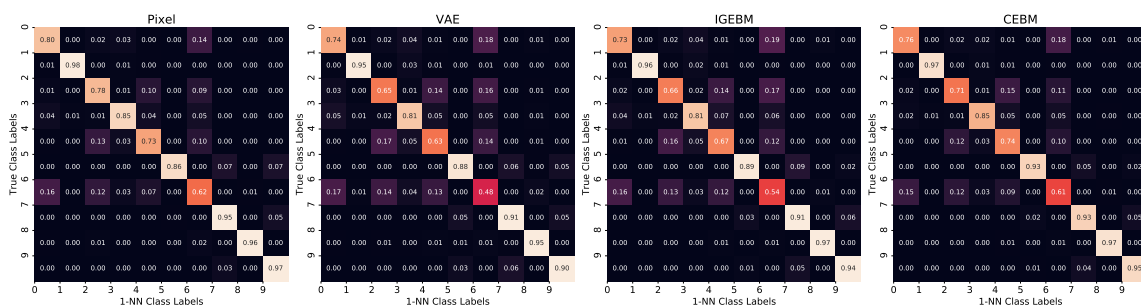


Figure 4: Fashion-MNIST

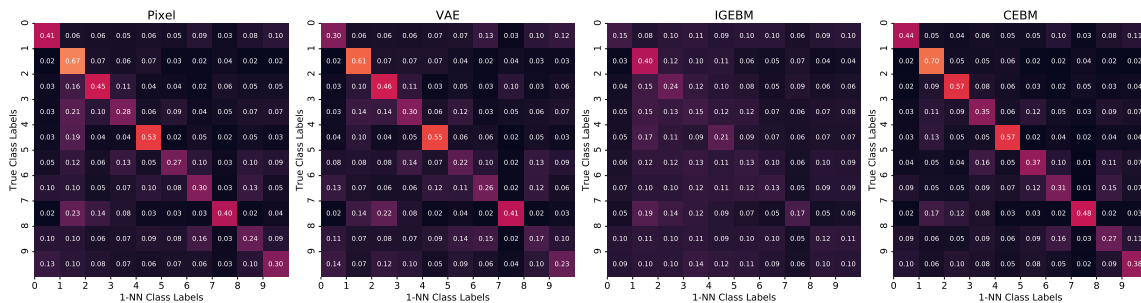


Figure 5: SVHN

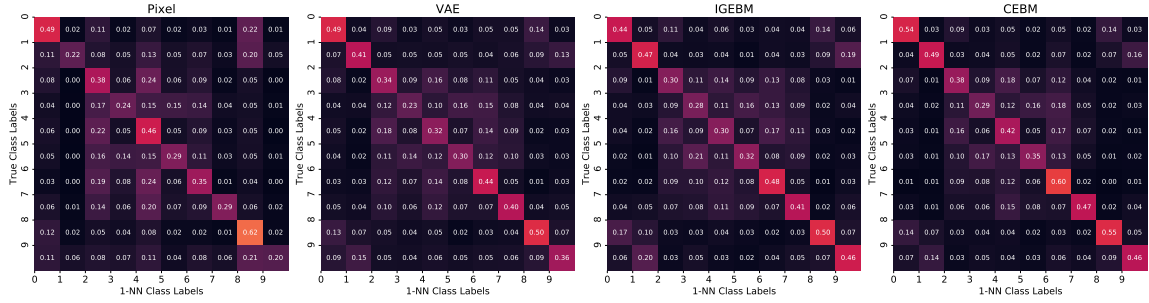


Figure 6: CIFAR10

Appendix G. Out-of-Distribution Detection Table

	Fashion-MNIST						CIFAR-10					
	$\log p_{\theta}(\mathbf{x})$			$\ \nabla_{\mathbf{x}} \log p_{\theta}(\mathbf{x})\ $			$\log p_{\theta}(\mathbf{x})$			$\ \nabla_{\mathbf{x}} \log p_{\theta}(\mathbf{x})\ $		
	MNIST	E-MNIST	C	MNIST	E-MNIST	C	SVHN	Texture	C	SVHN	Texture	C
VAE	50	39	9	61	57	1	42	58	41	38	51	37
IGE BM	35	36	90	78	82	96	45	31	64	33	17	62
CEBM	37	34	90	82	89	98	47	32	66	31	17	54
CEBMM	56	56	92	56	80	95	55	30	62	40	23	62

Table 5: AUROC scores in OOD Detection. We use $\log p_{\theta}(\mathbf{x})$ and $\|\nabla_{\mathbf{x}} \log p_{\theta}(\mathbf{x})\|$ as score functions. The left block shows results of the models trained on F-MNIST and tested on MNIST, E-MNIST, Constant (C); The right block shows results of the models trained on CIFAR-10 and tested on SVHN, Texture and Constant (C).