Is the Likelihood-based Deep Generative Modeling appropriate for Representation Learning?

Jakub M. Tomczak Copenhagen, October 12-13, 2021





What is data representation [1]:

"A good representation is (...) one that is useful as input to a supervised predictor."



What is data representation [1]:

"A good representation is (...) one that is useful as input to a supervised predictor."

$$\mathbf{z} = f_{\theta}(\mathbf{x})$$

 $y = f_{\varphi}(\mathbf{z})$

z is useful for the predictor = low error (?)

what about unsupervised learning?



What is data representation [1]:

"A good representation is (...) one that is useful as input to a supervised predictor."

$$\mathbf{z} = f_{\theta}(\mathbf{x})$$

 $y = f_{\varphi}(\mathbf{z})$

or, in the context of probabilistic modeling:

what about unsupervised learning?

"(...) a **good representation** is often one **that captures the posterior distribution of the underlying explanatory factors** for the observed input."



What is data representation [1]:

"A good representation is (...) one that is useful as input to a supervised predictor."

$$\mathbf{z} = f_{\theta}(\mathbf{x})$$

 $y = f_{\varphi}(\mathbf{z})$

or, in the context of probabilistic modeling:

what about unsupervised learning?

"(...) a **good representation** is often one **that captures the posterior distribution of the underlying explanatory factors** for the observed input."

$$q_{\phi}(\mathbf{z}|\mathbf{x})$$

q(**z**|**x**) is non-trivial (?), e.g., N(0, **1**)

- smoothness: $x \approx y$ implies $f(x) \approx f(y)$
- multiple explanatory factors
- a hierarchy of explanatory factors
- spatial/temporal coherence
- sparsity
- factors explain x and help to predict y
- geometrical properties (manifolds)
- generalizability (shared factors across tasks)
- simplicity of factor dependencies



Implic	it (model inductive biases)		<i>Explicit</i> (objective inductive biases)
• smooth	ness: $x \approx y$ implies $f(x) \approx f(y)$	•	smoothness: $x \approx y$ implies $f(x) \approx f(y)$
• multiple	e explanatory factors	•	sparsity
• a hierar	chy of explanatory factors	•	geometrical properties (manifolds)
• geomet	rical properties (manifolds)	•	(?) spatial/temporal coherence
• factors	explain x and help to predict y	•	(?) generalizability (shared factors
• (?) spat	al/temporal coherence		across tasks)
• (?) simp	licity of factor dependencies		



What makes a representation good (guidelines) [1]:

Implicit (model inductive biases)	<i>Explicit</i> (objective inductive biases)
smoothness: $x \approx y$ implies $f(x) \approx f(y)$	• smoothness: $x \approx y$ implies $f(x) \approx f(y)$
multiple explanatory factors	 sparsity
a hierarchy of explanatory factors	 geometrical properties (manifolds)
geometrical properties (manifolds)	• (?) spatial/temporal coherence
factors explain x and help to predict y	• (?) generalizability (shared factors
(?) spatial/temporal coherence	across tasks)
(?) simplicity of factor dependencies	
	<pre>Implicit (model inductive biases) smoothness: x ≈ y implies f(x) ≈ f(y) multiple explanatory factors a hierarchy of explanatory factors geometrical properties (manifolds) factors explain x and help to predict y (?) spatial/temporal coherence (?) simplicity of factor dependencies</pre>

What about compression?





smoothness: x		
multiple explaies is there a single fra	Is there a single framework for that?	
a hierarchy of Probabilistic modeli	ng +	











After [2]:

Generative modeling: formulating a model of the joint distribution, $p_{\vartheta}(\mathbf{x}, y)$

Discriminative modeling: formulating a model of the conditional, $p_{\vartheta}(y|\mathbf{x})$



15 [2] Jebara, T. (2012). Machine learning: discriminative and generative (Vol. 755). Springer Science & Business Media.

After [2]:

Generative modeling: formulating a model of the joint distribution, $p_{\vartheta}(\mathbf{x}, y)$

Discriminative modeling: formulating a model of the conditional, $p_{\vartheta}(y|\mathbf{x})$

Why do we need **generative modeling**?



16 [2] Jebara, T. (2012). Machine learning: discriminative and generative (Vol. 755). Springer Science & Business Media.

After [2]:

Generative modeling: formulating a model of the joint distribution, $p_{\vartheta}(\mathbf{x}, y)$

Discriminative modeling: formulating a model of the conditional, $p_{\vartheta}(y|\mathbf{x})$





After [2]:

Generative modeling: formulating a model of the joint distribution, $p_{\vartheta}(\mathbf{x}, y)$

Discriminative modeling: formulating a model of the conditional, $p_{\vartheta}(y|\mathbf{x})$



p(blue|x) is high
= certain decision!



After [2]:

Generative modeling: formulating a model of the joint distribution, $p_{\vartheta}(\mathbf{x}, y)$

Discriminative modeling: formulating a model of the conditional, $p_{\vartheta}(y|\mathbf{x})$



After [2]:

Generative modeling: formulating a model of the joint distribution, $p_{\vartheta}(\mathbf{x}, y)$

Discriminative modeling: formulating a model of the conditional, $p_{\vartheta}(y|\mathbf{x})$



Considering the joint could be crucial for decision making!

We consider the joint distribution factorized as follows:

$$p_artheta(\mathbf{x},y) = p_artheta(y|\mathbf{x}) \; p_artheta(\mathbf{x})$$

Then the log-likelihood function (assuming *iid*) takes the following form:

$$\ln p_{\vartheta}(\mathcal{D}) = \sum_{n} \ln p_{\vartheta}(y_n | \mathbf{x}_n) + \sum_{n} \ln p_{\vartheta}(\mathbf{x}_n)$$



21

We consider the joint distribution factorized as follows:

$$p_artheta(\mathbf{x},y) = p_artheta(y|\mathbf{x}) \; p_artheta(\mathbf{x})$$

Then the log-likelihood function (assuming *iid*) takes the following form:

$$\ln p_{\vartheta}(\mathcal{D}) = \sum_{n} \ln p_{\vartheta}(y_n | \mathbf{x}_n) + \sum_{n} \ln p_{\vartheta}(\mathbf{x}_n)$$

Parameterized by deep neural nets

Various models:

- \rightarrow Deep Latent Variable Models (e.g.,
- VAEs, diffusion models)
- \rightarrow Autoregressive Models (ARMs)
- \rightarrow Energy-based Models (EBMs) VU

We consider the joint distribution factorized as follows:

$$p_artheta(\mathbf{x},y) = p_artheta(y|\mathbf{x}) \; p_artheta(\mathbf{x})$$

Then the log-likelihood function (assuming *iid*) takes the following form:

$$\ln p_{\vartheta}(\mathcal{D}) = \sum_{n} \ln p_{\vartheta}(y_{n} | \mathbf{x}_{n}) + \sum_{n} \ln p_{\vartheta}(\mathbf{x}_{n})$$
We can use a shared parameterization.
$$p(y | \mathbf{x}) \quad p(\mathbf{x})$$
Neural Net
Neural Net
Neural Net



23

We consider the joint distribution factorized as follows:

$$p_artheta(\mathbf{x},y) = p_artheta(y|\mathbf{x}) \; p_artheta(\mathbf{x})$$

Then the log-likelihood function (assuming *iid*) takes the following form:

$$\ln p_{\vartheta}(\mathcal{D}) = \sum_{n} \ln p_{\vartheta}(y_n | \mathbf{x}_n) + \sum_{n} \ln p_{\vartheta}(\mathbf{x}_n)$$

1. Are there any problems with this objective?

2. Is this objective appropriate for representation learning?



The likelihood function



Let us take a look at the objective again (the log-likelihood function):

$$\ln p_{\vartheta}(\mathcal{D}) = \sum_{n} \ln p_{\vartheta}(y_n | \mathbf{x}_n) + \sum_{n} \ln p_{\vartheta}(\mathbf{x}_n)$$



Let us take a look at the objective again (the log-likelihood function):

$$\ln p_{\vartheta}(\mathcal{D}) = \sum_{n} \ln p_{\vartheta}(y_n | \mathbf{x}_n) + \sum_{n} \ln p_{\vartheta}(\mathbf{x}_n)$$

The conditional log-likelihood function is typically fine!



Let us take a look at the objective again (the log-likelihood function):

$$\ln p_{\vartheta}(\mathcal{D}) = \sum_{n} \ln p_{\vartheta}(y_n | \mathbf{x}_n) + \sum_{n} \ln p_{\vartheta}(\mathbf{x}_n)$$

The conditional log-likelihood function is typically **fine**!

A potential danger with the marginal log-likelihood function:



Let us take a look at the objective again (the log-likelihood function):

$$\ln p_{\vartheta}(\mathcal{D}) = \sum_{n} \ln p_{\vartheta}(y_n | \mathbf{x}_n) + \sum_{n} \ln p_{\vartheta}(\mathbf{x}_n)$$

The conditional log-likelihood function is typically **fine**!

A potential danger with the marginal log-likelihood function:

$$\rightarrow$$
 Consider a latent variable model: $p_{\vartheta}(\mathbf{x}) = \int p_{\vartheta}(\mathbf{x}|\mathbf{z}) p_{\vartheta}(\mathbf{z}) d\mathbf{z}$



Let us take a look at the objective again (the log-likelihood function):

$$\ln p_{\vartheta}(\mathcal{D}) = \sum_{n} \ln p_{\vartheta}(y_n | \mathbf{x}_n) + \sum_{n} \ln p_{\vartheta}(\mathbf{x}_n)$$

The conditional log-likelihood function is typically **fine**!

A potential danger with the marginal log-likelihood function:

 \rightarrow Consider a latent variable model: $p_{\vartheta}(\mathbf{x}) = \int p_{\vartheta}(\mathbf{x}|\mathbf{z}) \ p_{\vartheta}(\mathbf{z}) \ d\mathbf{z}$

 \rightarrow If we take Gaussians, then we get an **unbounded likelihood** [3]!

Solution: add a constant diagonal matrix to the decoder's covariance matrix.



Let us take a look at the objective again (the log-likelihood function):

$$\ln p_{\vartheta}(\mathcal{D}) = \sum_{n} \ln p_{\vartheta}(y_n | \mathbf{x}_n) + \sum_{n} \ln p_{\vartheta}(\mathbf{x}_n)$$

The conditional log-likelihood function is typically **fine**!

A potential danger with the marginal log-likelihood function:

- \rightarrow Consider a latent variable model: $p_{\vartheta}(\mathbf{x}) = \int p_{\vartheta}(\mathbf{x}|\mathbf{z}) p_{\vartheta}(\mathbf{z}) d\mathbf{z}$
- \rightarrow If we take Gaussians, then we get an **unbounded likelihood** [3]!
- \rightarrow If we take Bernoullis (or Categoricals), it's **fine** [3]!

Is then the likelihood function well-suited for representation learning?



Is then the likelihood function well-suited for representation learning?

$$\begin{split} KL[p_{data}(\mathbf{x})||p_{\theta}(\mathbf{x})] &= -\mathbb{H}[p_{data}(\mathbf{x})] + \mathbb{CE}[p_{data}(\mathbf{x})||p_{\theta}(\mathbf{x})] \\ &= const + \mathbb{CE}[p_{data}(\mathbf{x})||p_{\theta}(\mathbf{x})] \\ &= const - \int p_{data}(\mathbf{x})\ln p_{\theta}(\mathbf{x}) \, \mathrm{d}\mathbf{x} \\ &= const - \frac{1}{N}\sum_{n=1}^{N}\ln p_{\theta}(\mathbf{x}_{n}) \end{split}$$

The likelihood-based approach is equivalent to calculating a match between a model and the empirical distribution according to the KL-divergence.



Is then the likelihood function well-suited for representation learning?

<i>Implicit</i> (model inductive biases)	<i>Explicit</i> (objective inductive biases)
	• smoothness: $x \approx y$ implies $f(x) \approx f(y)$
	• sparsity
	 geometrical properties (manifolds)
	• (?) spatial/temporal coherence
	• (?) generalizability (shared factors
	across tasks)

The likelihood function alone does <u>not</u> help representation learning *explicitly*.



Is then the likelihood function well-suited for representation learning? Three scenarios:





Is then the likelihood function well-suited for representation learning? Three scenarios:



The likelihood function has no incentive to learn useful representations.




The likelihood function for representation learning

Is then the likelihood function well-suited for representation learning? Three scenarios:



The likelihood function has no incentive to learn *useful* representations.

But is it bad? Maybe we should focus on the parameterization instead!



The likelihood function for representation learning

Is then the likelihood function well-suited for representation learning?

Implicit (model inductive biases)	<i>Explicit</i> (objective inductive biases)
• smoothness: $x \approx y$ implies $f(x) \approx f(y)$	
 multiple explanatory factors 	
• a hierarchy of explanatory factors	
 geometrical properties (manifolds) 	
• factors explain x and help to predict y	
• (?) spatial/temporal coherence	
• (?) simplicity of factor dependencies	

The likelihood function is a *general* objective and the parameterization is key?

How to model the marginal distribution? What parameterization?



Classes of DGMs

Before we delve into the parameterizations, let us recap the DGMs:

- Latent Variables Models: $p(\mathbf{x}) = \int p(\mathbf{x}|\mathbf{z}) \ p(\mathbf{z}) \ d\mathbf{z}$
- Flow-based models: $p(\mathbf{x}) = p(\mathbf{z} = f(\mathbf{x}))|\mathbf{J}_{f(\mathbf{x})}|$

• Autoregressive models:
$$p(\mathbf{x}) = p(x_0) \prod_{i=1}^{D} p(x_i | \mathbf{x}_{< i})$$

• Energy-based models: $p(\mathbf{x}) = \frac{\exp\{-E(\mathbf{x})\}}{Z}$



Classes of DGMs

Before we delve into the parameterizations, let us recap the DGMs:

- Latent Variables Models: $p(\mathbf{x}) = \int p(\mathbf{x}|\mathbf{z}) \, p(\mathbf{z}) \, d\mathbf{z}$
- Flow-based models: $p(\mathbf{x}) = p(\mathbf{z} = f(\mathbf{x}))|\mathbf{J}_{f(\mathbf{x})}|$
- Latent variables = representation

- Autoregressive models: $p(\mathbf{x}) = p(x_0) \prod p(x_i | \mathbf{x}_{< i})$
- Energy-based models. $p(\mathbf{x}) = \frac{\exp\{-E(\mathbf{x})\}}{Z}$

What is a *representation* here?



Importance of parameterization in Flows

In flows, we need either powerful **invertible** neural networks, e.g.:



Importance of parameterization in Flows

In flows, we need either powerful invertible neural networks, e.g.:

 \rightarrow ResNets [5] or DenseNets [6] with constraints (less than 1-Lipschitz).



[5] Chen, R. T., Behrmann, J., Duvenaud, D. K., & Jacobsen, J. H. (2019). Residual Flows for Invertible Generative Modeling. NeurIPS 2019
 [6] Perugachi-Diaz, Y., Tomczak, J. M., & Bhulai, S. (2021). Invertible densenets with concatenated lipswish. NeurIPS 2021



Importance of parameterization in Flows

In flows, we need either powerful **invertible** neural networks, e.g.:

 \rightarrow ResNets [5] or DenseNets [6] with constraints (less than 1-Lipschitz).

or/and invertible structure with easily computable Jacobian, e.g.:

 \rightarrow Householder flows [7], Sylvester flows [8]

 \rightarrow Convolution Exponential + Sylvester flows [9]

[5] Chen, R. T., Behrmann, J., Duvenaud, D. K., & Jacobsen, J. H. (2019). Residual Flows for Invertible Generative Modeling. NeurIPS 2019
[6] Perugachi-Diaz, Y., Tomczak, J. M., & Bhulai, S. (2021). Invertible densenets with concatenated lipswish. NeurIPS 2021
[7] Tomczak, J. M., & Welling, M. (2016). Improving variational auto-encoders using householder flow. Bayesian Deep Learing @ NeurIPS 2016
[8] van den Berg, R., Hasenclever, L., Tomczak, J. M., & Welling, M. Sylvester Normalizing Flows for Variational Inference. UAI 2018
[9] Hoogeboom, E., Garcia Satorras, V., Tomczak, J., & Welling, M. (2020). The Convolution Exponential and Generalized Sylvester Flows. NeurIPS 2020



We can use flows to formulate a joint distribution [10]:

- The marginal distribution is a flow-based model.
- The mapping from **x** to **z** is shared with a predictor.
- We need to change the objective function:

$$\ell_{\lambda}(\mathbf{x}, y) = \ln p(y|\mathbf{x}) + \lambda \ln p(\mathbf{x})$$





We can use flows to formulate a joint distribution [10]:

- The marginal distribution is a flow-based model.
- The mapping from **x** to **z** is shared with a predictor.
- We need to change the objective function:

$$\ell_{\lambda}(\mathbf{x}, y) = \ln p(y|\mathbf{x}) + \lambda \ln p(\mathbf{x})$$

The gradient flow is uneven from the conditional and the marginal.





We see that λ is crucial [6]:

	$\lambda = 0$	$\lambda = \frac{1}{D}$		$\lambda = 1$	
Model \Evaluation	Acc \uparrow	Acc \uparrow	bpd \downarrow	Acc \uparrow	bpd \downarrow
Coupling	89.77%	87.58%	4.30	67.62%	3.54
$+1 \times 1$ conv	90.82%	87.96%	4.09	67.38%	3.47
Residual Blocks (full)	91.78%	90.47%	3.62	70.32%	3.39
Dense Blocks (full)	92.40%	90.79%	3.49	75.67%	3.31



(a) Hybrid model trained with $\lambda = \frac{1}{D}$



(b) Hybrid model trained with $\lambda=1$



We see that λ is crucial [6]:

	$\lambda = 0$	$\lambda = \frac{1}{D}$		$\lambda = 1$	
Model \Evaluation	Acc \uparrow	Acc \uparrow	bpd \downarrow	Acc \uparrow	bpd \downarrow
Coupling	89.77%	87.58%	4.30	67.62%	3.54
$+1 \times 1$ conv	90.82%	87.96%	4.09	67.38%	3.47
Residual Blocks (full)	91.78%	90.47%	3.62	70.32%	3.39
Dense Blocks (full)	92.40%	90.79%	3.49	75.67%	3.31

 $\ell_{\lambda}(\mathbf{x}, y) = \ln p(y|\mathbf{x}) + \lambda \ln p(\mathbf{x})$

This is not the likelihood function anymore...



(a) Hybrid model trained with $\lambda = \frac{1}{D}$



(b) Hybrid model trained with $\lambda=1$

Variational Auto-Encoders

• Variational Auto-Encoders (VAEs) constitute a broad class of models.



Variational Auto-Encoders

- Variational Auto-Encoders (VAEs) constitute a broad class of models.
- But they fail to learn a *representation* [11].



[11] Hübotter, J. F., Lanillos, P., & Tomczak, J. M. (2021). Training Deep Spiking Auto-encoders without Bursting or Dying Neurons through 50 Regularization. arXiv preprint arXiv:2109.11045.

Variational Auto-Encoders

- Variational Auto-Encoders (VAEs) constitute a broad class of models.
- But they fail to learn a *representation* [11].
- It is possible to introduce geometrical perspective to VAEs:
- → Hyperspherical VAEs [12], Riemaniann manifolds [13], Lie groups [14]
- Hierarchical latent structure is crucial for learning powerful VAEs.

[11] Hübotter, J. F., Lanillos, P., & Tomczak, J. M. (2021). Training Deep Spiking Auto-encoders without Bursting or Dying Neurons through Regularization. arXiv preprint arXiv:2109.11045.
[12] Davidson, T. R., Falorsi, L., De Cao, N., Kipf, T., & Tomczak, J. M. (2018). Hyperspherical variational auto-encoders. UAI 2018
[13] Arvanitidis, G., Hauberg, S., & Schölkopf, B. (2021). Geometrically Enriched Latent Spaces. AISTATS 2021
[14] Falorsi, L., de Haan, P., Davidson, T. R., & Forré, P. (2019). Reparameterizing distributions on lie groups. AISTATS 2019



A proper **parameterization** of hierarchical VAEs is important to get SOTA. So called top-down VAEs:

→ ResNet VAE [15], Ladder VAE [16], BIVA [17], NVAE [18], vdVAE [19]

 $p(\mathbf{x}|\mathbf{z}_1)p(\mathbf{z}_1|\mathbf{z}_2)p(\mathbf{z}_2)$ $q(\mathbf{z}_1|\mathbf{z}_2,\mathbf{x})q(\mathbf{z}_2|\mathbf{x})$



[15] Kingma, D. P., Salimans, T., Jozefowicz, R., Chen, X., Sutskever, I., & Welling, M. (2016). Improved variational inference with inverse autoregressive flow. NeurIPS 2016 [16] Sønderby, C. K., Raiko, T., Maaløe, L., Sønderby, S. K., & Winther, O. (2016). Ladder variational autoencoders. NeurIPS 2016 [17] Maaløe, L., Fraccaro, M., Liévin, V., & Winther, O. (2019). BIVA: A Very Deep Hierarchy of Latent Variables for Generative Modeling. NeurIPS 2019 [18] Vahdat, A., & Kautz, J. (2020). Nvae: A deep hierarchical variational autoencoder. NeurIPS 2020 [19] Child, R. (2020). Very Deep VAEs Generalize Autoregressive Models and Can Outperform Them on Images. ICLR 2020



VAEs for joint modeling

However, still, using VAEs for modeling the marginal and sharing the encoder with the predictor requires applying a fudge factor.

The objective can take the following form [20, 21]:

$$\ell_{\alpha}(\mathbf{x}, y) = \alpha \mathbb{E}_{q(\mathbf{z}|\mathbf{x})} \left[\ln p(y|\mathbf{z}) \right] + ELBO(\mathbf{x})$$

The encoder is shared with the classifier.



VAEs for joint modeling

However, still, using VAEs for modeling the marginal and sharing the encoder with the predictor requires applying a fudge factor.

The objective can take the following form [20, 21]:

$$\ell_{\alpha}(\mathbf{x}, y) = \alpha \mathbb{E}_{q(\mathbf{z}|\mathbf{x})} \left[\ln p(y|\mathbf{z}) \right] + ELBO(\mathbf{x})$$

The encoder is shared with the classifier.

This is not the likelihood function anymore...

[20] Kingma, D. P., Mohamed, S., Rezende, D. J., & Welling, M. (2014). Semi-supervised learning with deep generative models. NeurIPS 2014
 [21] Ilse, M., Tomczak, J. M., Louizos, C., & Welling, M. (2020). DIVA: Domain invariant variational autoencoders. MIDL 2020



It is unclear for ARMs and EBMs what a *representation* is.



It is unclear for ARMs and EBMs what a *representation* is.

However, for EBMs we have [22]:

$$\ln p_{\theta}(\mathbf{x}, y) = \ln p_{\theta}(y|\mathbf{x}) + \ln p_{\theta}(\mathbf{x})$$

$$= \ln \frac{\exp\{NN_{\theta}(\mathbf{x})[y]\}}{\sum_{y} \exp\{NN_{\theta}(\mathbf{x})[y]\}} + \ln \frac{\sum_{y} \exp\{NN_{\theta}(\mathbf{x})[y]\}}{Z_{\theta}}$$

$$= \ln \operatorname{softmax}\{NN_{\theta}(\mathbf{x})[y]\} + \left(\operatorname{LSE}_{y}\{NN_{\theta}(\mathbf{x})[y]\} - \ln Z_{\theta}\right)$$

$$NN_{\theta}(\mathbf{x})$$

[22] Grathwohl, W., Wang, K. C., Jacobsen, J. H., Duvenaud, D., Norouzi, M., & Swersky, K. (2020). Your Classifier is Secretly an Energy Based Model and You Should Treat it Like One. ICLR 2020



 $\tilde{p}_{\theta}(\mathbf{x})$

LogSumExp₂

 $\bigcirc NN_{\theta}(\mathbf{x})[y]$

 $p_{\theta}(y|\mathbf{x})$

 (\mathbf{x})

It is unclear for ARMs and EBMs what a *representation* is.

However, for EBMs we have [22]:

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





[22] Grathwohl, W., Wang, K. C., Jacobsen, J. H., Duvenaud, D., Norouzi, M., & Swersky, K. (2020). Your Classifier is Secretly an Energy Based Model and You Should Treat it Like One. ICLR 2020

57

It is unclear for ARMs and EBMs what a *representation* is.

However, for EBMs we have [22]:

$$\ln p_{\theta}(\mathbf{x}, y) = \ln p_{\theta}(y|\mathbf{x}) + \ln p_{\theta}(\mathbf{x})$$
$$= \ln \frac{\exp\{NN_{\theta}(\mathbf{x})[y]\}}{\sum_{y} \exp\{NN_{\theta}(\mathbf{x})[y]\}} + \ln \frac{\sum_{y} \exp\{NN_{\theta}(\mathbf{x})[y]\}}{Z_{\theta}}$$
$$= \ln \operatorname{softmax}\{NN_{\theta}(\mathbf{x})[y]\} + \left(\operatorname{LSE}_{y}\{NN_{\theta}(\mathbf{x})[y]\} - \ln Z_{\theta}\right)$$

A single, shared neural net.

The objective is the likelihood function!

[22] Grathwohl, W., Wang, K. C., Jacobsen, J. H., Duvenaud, D., Norouzi, M., & Swersky, K. (2020). Your Classifier is Secretly an Energy Based Model and You Should Treat it Like One. ICLR 2020





It is unclear for ARMs and EBMs what a *representation* is.

However, for EBMs we have [22]:

$$\begin{aligned} \ln p_{\theta}(\mathbf{x}, y) &= \ln p_{\theta}(y | \mathbf{x}) + \ln p_{\theta}(\mathbf{x}) \\ &= \ln \frac{\exp\{NN_{\theta}(\mathbf{x})[y]\}}{\sum_{y} \exp\{NN_{\theta}(\mathbf{x})[y]\}} + \ln \frac{\sum_{y} \exp\{NN_{\theta}(\mathbf{x})[y]\}}{Z_{\theta}} \\ &= \ln \operatorname{softmax}\{NN_{\theta}(\mathbf{x})[y]\} + \left(\operatorname{LSE}_{y}\{NN_{\theta}(\mathbf{x})[y]\} - \ln Z_{\theta}\right) \end{aligned}$$

A single, shared neural net.

The objective is the likelihood function!

But what is the representation here?





It is unclear for ARMs and EBMs what a *representation* is.

However, for EBMs we have [22]:

$$\ln p_{\theta}(\mathbf{x}, y) = \ln p_{\theta}(y|\mathbf{x}) + \ln p_{\theta}(\mathbf{x})$$

$$= \ln \frac{\exp\{NN_{\theta}(\mathbf{x})[y]\}}{\sum_{y} \exp\{NN_{\theta}(\mathbf{x})[y]\}} + \ln \frac{\sum_{y} \exp\{NN_{\theta}(\mathbf{x})[y]\}}{Z_{\theta}}$$

$$= \ln \operatorname{softmax}\{NN_{\theta}(\mathbf{x})[y]\} + \left(\operatorname{LSE}_{y}\{NN_{\theta}(\mathbf{x})[y]\} - \ln Z_{\theta}\right)$$

$$NN_{\theta}(\mathbf{x})$$

$$A \text{ single, shared neural net.}$$

$$Iogits?$$

$$NN_{\theta}(\mathbf{x})$$

$$How to obtain a representation?$$

 $p_{\theta}(y|\mathbf{x})$

 $\tilde{p}_{\theta}(\mathbf{x})$

[22] Grathwohl, W., Wang, K. C., Jacobsen, J. H., Duvenaud, D., Norouzi, M., & Swersky, K. (2020). Your Classifier is Secretly an Energy Based Model and You Should Treat it Like One. ICLR 2020

The likelihood function:

- Potential issues with mixture models (unboundness).
- Fine for discrete distributions.



The likelihood function:

- Potential issues with mixture models (unboundness).
- Fine for discrete distributions.

Marginal models:

- For ARMs and EBMs: Unclear what a *representation* is.
- For VAEs and flows: We don't know how to use the likelihood function.



The likelihood function:

- Potential issues with mixture models (unboundness).
- Fine for discrete distributions.

Marginal models:

Do we really need one?

- For ARMs and EBMs: Unclear what a *representation* is.
- For VAEs and flows: We don't know how to use the likelihood function.

Where is the problem?



Do we need a notion of *representation learning*?



What is a *representation*?

Learning a joint distribution raises questions:

- Is there something wrong with the likelihood function? Or rather:
- How should we parameterize distributions?

Interestingly, EBMs don't need to define an explicit representation and they can use the likelihood function!



What is a *representation*?

Learning a joint distribution raises questions:

- Is there something wrong with the likelihood function? Or rather:
- How should we parameterize distributions?

Interestingly, EBMs don't need to define an explicit representation and they can use the likelihood function!

So do we really need a *representation*?

Or maybe we misinterpret the *representation*?



Shannon's source coding theorem:

The length of a message representing some data is proportional to the

entropy of this data.



Shannon's source coding theorem:

The length of a message representing some data is proportional to the

entropy of this data.

In other words:

Maximizing the model log-likelihood is equivalent to minimizing the expected number of bits required per message, if the encoder is optimal.

So we want to get $p_{real}(\mathbf{x}) \approx p_{\theta}(\mathbf{x})$.



Shannon's theorem could be read as follows:



Shannon's theorem could be read as follows:

 \rightarrow A model maximizing the likelihood function is **useful for compression**.



Shannon's theorem could be read as follows:

 \rightarrow A model maximizing the likelihood function is **useful for compression**.

 \rightarrow There is no need to define a separate variable to send a message.

A good model $p_{\theta}(\mathbf{x})$ is all we need!

Shannon's theorem could be read as follows:

 \rightarrow A model maximizing the likelihood function is **useful for compression**.

 \rightarrow There is no need to define a separate variable to send a message.

A good model $p_{\theta}(\mathbf{x})$ is all we need!

 \rightarrow A message is discrete in nature (i.e., symbols and bits).


Compression

Shannon's theorem could be read as follows:

 \rightarrow A model maximizing the likelihood function is **useful for compression**.

 \rightarrow There is no need to define a separate variable to send a message.

A good model $p_{\theta}(\mathbf{x})$ is all we need!

 \rightarrow A message is discrete in nature (i.e., symbols and bits).

In information theory we have a "natural" notion of compressed representation of data. Do we need more?



Compression

Shannon's theorem could be read as follows:

 \rightarrow A model maximizing the likelihood function is **useful for compression**.

 \rightarrow There is no need to define a separate variable to send a message.

A good model $p_{\theta}(\mathbf{x})$ is all we need!

 \rightarrow A message is discrete in nature (i.e., symbols and bits).

In information theory we have a "natural" notion of compressed representation of data. Do we need more?

But then, how to use compression in the context of prediction?





- 1. Is the likelihood function enough?
 - a. If NO, then what priors/regularizers are important? (*Explicit*)
 - b. If YES, what parameterizations are important? (Implicit)



- 1. Is the likelihood function enough?
 - a. If NO, then what priors/regularizers are important? (*Explicit*)
 - b. If YES, what parameterizations are important? (Implicit)



- 1. Is the likelihood function enough?
 - a. If NO, then what priors/regularizers are important? (*Explicit*)
 - b. If YES, what parameterizations are important? (Implicit)

2. How to obtain a *representation* from EBMs or ARMs?



- 1. Is the likelihood function enough?
 - a. If NO, then what priors/regularizers are important? (*Explicit*)
 - b. If YES, what parameterizations are important? (*Implicit*)

2. How to obtain a *representation* from EBMs or ARMs?

3. Is the notion of *representation learning* even necessary? Is it a different label for compression?



Thank you!

