Evaluating Deep Generative Models A p(reality) check

Joey Bose





Generative Modelling



Image Super Resolution

"Creating noise from data is easy; creating data from noise is generative modeling." Song et. Al 2020



Text to Speech

Drug Discovery

Types of Generative Models

• Fully-observed models

- Latent Variable Models
 - Prescribed Models: Likelihood + noise
 - Implict Models (Likelihood-free)



Score Matching

Fig Credit: Song et. Al 2021

Data

KLD

MMD

JSD

Fig Credit: Theis et. al 2016

Data

KLD

Avoids assigning extremely small probability to any data point but assigns a lot of probability mass to non-data regions.

MMD

Fig Credit: Theis et. al 2016

JSD

Data

MMD

Fig Credit: Theis et. al 2016

Data

KLD

Different Evaluation Mesures have Different Tradeoffs!

MMD

JSD

Fig Credit: Theis et. al 2016

Key question: What is the task you care about?

- Density Estimation
- Sampling/Generation
- Semi-supervised Learning
- Learning useful representations
- Anomaly Detection
- Hybrid Modeling
- Calibration

Key question: What is the task you care about?

- Density Estimation
- Sampling/Generation
- Semi-supervised Learning
- Learning useful representations
- Anomaly Detection
- Hybrid Modeling
- Calibration

Most research papers do this

Key question: What is the task you care about?

- Density Estimation
- Sampling/Generation
- Semi-supervised Learning
- Learning useful representations
- Anomaly Detection
- Hybrid Modeling
- Calibration

Most research papers do this Others claim their models do a good job of this

Key question: What is the task you care about?

- Density Estimation
- Sampling/Generation
- Semi-supervised Learning
- Learning useful representations
- Anomaly Detection
- Hybrid Modeling
- Calibration

Most research papers do this Others claim their models do a good job of this Hope we will magically do well on these.

Key question: What is the task you care about?

- Density Estimation
- Sampling/Generation
- Semi-supervised Learning
- Learning useful representations
- Anomaly Detection
- Hybrid Modeling
- Calibration

How much progress have we actually made on these tasks?

Most research papers do this Others claim their models do a good job of this Hope we will magically do well on these.

Explicit Likelihood Models for Anomaly Detection

So you have an explicit likelihood model and you want to do Anomaly Detection? Should be simple right?

Problem Definition

- Input space $\mathscr{X} \subseteq \mathbb{R}^D$
- Partitioning $\mathscr{X}_{in} \cup \mathscr{X}_{out} = \mathscr{X}$
- $\mathscr{X}_{in} \cap \mathscr{X}_{out} = \emptyset$

Pick \mathcal{X}_{in} such that it contains "the majority" of the mass e.g. $P_X^*(\mathcal{X}_{in}) = 1 - \alpha \in (0.5, 1)$

x

Fig Credit: Le Lan and Dinh 2021

Approaches for Anomaly Detection Density Scoring

Unlikely samples should have "low" likelihood —i.e. low $p_X^*(x)$ threshold $\lambda > 0$

A heuristic is to define inliers as points whose density is greater than a

Fig Credit: Le Lan and Dinh 2021

Approaches for Anomaly Detection Typicality Test

 $\leq \epsilon$

close to the average log-density of the distribution $\lambda > 0$. Loose definition, the typical set $A_{\epsilon}^{(N)}(p_{x}^{*}) \subset \mathscr{X}^{N}$ of points satisfy:

$$H(p^*x) + \frac{1}{N} \sum_{n=1}^{N} \log p_X^*(x^{(n)})$$

Inliers are part of typical set --i.e. points whose average log density is

Fig Credit: Le Lan and Dinh 2021

Approaches for Anomaly Detection

Assume we have trained the perfect and capacity: $p_X^{(\theta)} = p_X^*$

The observed data is a chosen representation of the real data and should remain invariant under an invertible map f (There is no loss of information).

$$C^*(x)$$

Perfect Classifier on X

Assume we have trained the perfect density model under infinite data

$$(c^* \circ f^{-1})(f(x))$$

Perfect Classifier on f(X)

Approaches for Anomaly Detection

Principle. In an infinite data and capacity setting, the result of an

 \mathscr{X}_{out} remains a low probability subset as $P_X(\mathcal{X}_{out}) = P_{f(X)}(f(\mathcal{X}_{out})) \quad \text{and} \quad \forall x \in \mathcal{X}, x \in \mathcal{X}_{out} \Longleftrightarrow f(x) \in f(\mathcal{X}_{out})$

 $p_{f(X)}^{*}(f(x)) = p_{X}^{*}(x) \left| \frac{\partial f}{\partial x^{T}}(x) \right|^{-1}$ \longrightarrow Change density under an invertible map

Fig Credit: Le Lan and Dinh 2021

- anomaly detection should be invariant to continuous reparametrization.

Anomaly Detection: Uniformization

Under weak assumptions one can map any distribution to an uniform one with an invertible map f^{KR} . $p^*_{f(KR)(X)} = 1$ is constant everywhere.

$$\forall d \in \{1, \cdots, D\}, f^{(KR)}(x) = CDF_{p^*x_{d|X_{< d}}}(x_d | x_{< d})$$

Anomaly Detection: Arbitrary Score Proposition 1(Le Lan and Dinh 2021). For any R.V. $X \sim p_X^*$ with p_X^* continuous strictly positive (with \mathscr{X} convex) and any measurable continuous map $s: \mathscr{X} \to \mathbb{R}^*_+$ bounded below by a strictly positive number, there exists a continuous bijection $f^{(s)}$ s.t. $x \in \mathcal{X}, p_{f^{(s)}(X)}(f^{(s)}(x)) = s(x)$

Fig Credit: Le Lan and Dinh 2021

Hybrid Modeling

$\log p(x, y) = \log p(y \mid x) + \log p(x)$

The discriminative part will assign a score to any input even if the point is actually not part of the data distribution. Can we balance this with the generative part? Can we use any pre-trained generative model to help here?

What if we want to use generative models for discriminative tasks?

Generative Part

Discriminative part

Hybrid Modeling: Naive Approach If the discriminative and generative parts were trained separately we have: **Generative Part**

$$\log p(x, y) = \log p_{\theta_1}(y \mid x) + \log p_{\theta_2}(x)$$

Operationally we can use separate networks

Discriminative part

Hybrid Modeling: Naive Approach If the discriminative and generative parts were trained separately we have: **Generative Part**

$$\log p(x, y) = \frac{\log p_{\theta_1}(y \mid x)}{\log p_{\theta_2}(x)} + \log p_{\theta_2}(x)$$

There is no information flow between x and y. There is no reason each NN treats x in the same way!

Discriminative part

Hybrid Modeling What if we share parameters

Generative Part

Hybrid Modeling: Parameter Sharing

If y is a binary label it needs only 1 bit while for a D-dim x vector we need at least as many bits.

 $\nabla_{\gamma} logp(x, y) = \nabla_{\gamma} \log y$

Let's also assume that the output of both discriminative and generative nets is 0.5. What is the

$\log p_{\theta_1,\gamma}(y \mid 0.5) = y \log 0.5 + (1 - y) \log 0.5$

$$p_{\theta_1,\gamma}(y \mid x) + \nabla_{\gamma} \log p_{\theta_2,\gamma}(x)$$

Hybrid Modeling: Parameter Sharing

$\log Bern(y \mid 0.5) = y \log 0.5 + (1 - y) \log 0.5$ $= -\log 2$

 $= -D \log 2$

Hybrid Modeling: How do we fix this?

- al., 2019)

Take home message: If this is your end downstream task it can't be treated as two separate problems. The generative modelling part must be tightly coupled with the end task!

 Convex combination of the two objectives (Bouchard & Triggs, 2004) • Up weight the generative part by a positive factor $\lambda \ge 0$ (Nalisnick et

• Use an invertible network and a learned flow prior (Chen et al., 2019)

Generative Models as Lossy Compression

Shannon's fundamental compression theorem states that we can compress a random variable $x \sim p(x)$ losslessly at $\mathcal{H}(x)$. That is a a rate close the entropy for optimal reconstruction.

What is the tradeoff now?

- good generative model should be able to "transmit" a coded input with

What about lossy compression? Suppose we allow a compression rate $R \leq \mathcal{H}(x)$ using a code z and have a lossy reconstruction $\hat{x} = f(z)$.

Generative Models as Lossy Compression

Given a distortion threshold D Shannon's rate distortion theorem states that the rate distortion function $\mathscr{R}(D)$ equals the minimum of the following:

$$\min \mathcal{I}(x,z) \quad s \cdot t$$

$$q(z|x)$$
The encoder distribution or
"approximate posterior"

"

 $\mathbb{E}_{q(x,z)}[d(x,f(z))] \leq D$

Joint distribution

Latent Variable modeling

$$\log p(\mathcal{D}) = \sum_{i=1}^{N} \log \mathbb{E}_{p(z)} [p(x_i | x_i)]$$

$$\log p(\mathcal{D}) \ge \sum_{i=1}^{N} \mathbb{E}_{q_i(z)} \left[\log \frac{p(x | x_i)}{q_i}\right]$$

$$\mathbb{E}_{q_i(z)} \left[\log \frac{p(x \mid z) p(z)}{q_i(z)} \right] = \mathbb{E}_{q_i(z)} \left[1 - \frac{p(x \mid z) p(z)}{q_i(z)} \right]$$

Lossy Compression in VAE's

We can modify the formalism of rate distortion theory to match the generative model formalism

 $\mathscr{I}(x,z) \leq \mathscr{I}(x,z) + KL(q(z) | | p(z))$ $= \mathbb{E}_{p_d(x)}[KL(q(z \mid x) \mid | p(z))]$ "approximate posterior"

 $\min \mathbb{E}_{p_d(x)}[KL(q(z \mid x) \mid | p(z \mid x) \mid x)]$ q(z|x)

Variational Rate Distortion functio

$$z))] \quad s.t.\mathbb{E}_{q(x,z)}[d(x,f(z))] \leq$$

on:
$$\mathscr{R}_p(D)$$

Lossy Compression in VAE's

We can change the constrained optimization using the method of longrange multipliers:

$\min_{q(z|x)} \mathbb{E}_{p_d(x)}[KL(q(z|x) | | p(z)$

 $\min_{q(z|x)} \mathbb{E}_{p_d(x)} [KL(q(z|x) | | p(z)$

$$))] \quad s.t. \mathbb{E}_{q(x,z)}[d(x,f(z))] \leq \\ \downarrow \\))] + \beta \mathbb{E}_{q(x,z)}[d(x,f(z))] \leq D$$

Independent optimization problems for each x

Rate Distortion Curves

• $\mathscr{R}_p(D)$ is an upper bound for any prior on $\mathscr{R}(D)$ p(z)

$\mathscr{R}(D) = \min \mathscr{R}(D)$ which means for any β there's an optimal prior

Fig Credit: Huang et. al 2020

Rate Distortion Curves

structions ($\beta = 5$). (c) Variational rate distortion curves for each of these three models.

- $p_1(x), p_2(x)$ and $p_3(x)$ are 2-d gen models with 1-d latent code. Conditional likelihoods are the grey curves and coloured dots are
- prior samples
- $p_1(x)$ and $p_2(x)$ have the same prior but different decoder

Figure 2. (a) Prior samples, or equivalently, reconstructions with $\beta = 0$. (b) High-rate recon-

Rate Distortion Curves: GANs

- this.
- network make better use of the information in code space

Increasing the code size has the effect of extending the curve leftward. High-rate regime is effectively measuring reconstruction ability and a larger code size helps with

Increasing the depth pushes the curves down and to the left. Capacity helps the

What about Disentanglement?

 \mathcal{X}

informative factors of variations in the data (Bengio et. Al 2013)

is sampled from a distribution P(z). z corresponds to semantically

Principle. A disentangled representation should separate the distinct,

2-step Generative Process: First, a multivariate latent random variable zmeaningful factors of variation of the observations. Then, in a second step, the observation x is sampled from the conditional distribution $P(x \mid z)$.

What about Disentanglement?

Theorem 1. For d > 1, let $z \sim P$ denote any distribution which admits a density $p(\mathbf{z}) = \prod_{i=1}^{d} p(\mathbf{z}_i)$. Then, there exists an infinite family of bijective functions $f : \operatorname{supp}(\mathbf{z}) \rightarrow$ supp(z) such that $\frac{\partial f_i(u)}{\partial u_i} \neq 0$ almost everywhere for all i and j (i.e., z and f(z) are completely entangled) and $P(\mathbf{z} \leq \mathbf{u}) = P(f(\mathbf{z}) \leq \mathbf{u})$ for all $\mathbf{u} \in \operatorname{supp}(\mathbf{z})$ (i.e., they have the same marginal distribution).

Unsupervised Disentanglement is impossible (Locatello et. Al 2019)!

Strong statements but there are caveats.

Interpreting the Impossibility result

Theorem 1. For d > 1, let $z \sim P$ denote any distribution which admits a density $p(\mathbf{z}) = \prod_{i=1}^{d} p(\mathbf{z}_i)$. Then, there exists an infinite family of bijective functions $f : \operatorname{supp}(\mathbf{z}) \rightarrow$ $\operatorname{supp}(\mathbf{z})$ such that $\frac{\partial f_i(\mathbf{u})}{\partial u_i} \neq 0$ almost everywhere for all i and j (i.e., z and f(z) are completely entangled) and $P(\mathbf{z} \leq \mathbf{u}) = P(f(\mathbf{z}) \leq \mathbf{u})$ for all $\mathbf{u} \in \operatorname{supp}(\mathbf{z})$ (i.e., they have the same marginal distribution).

Suppose you have disentanglement method that gives a representation r(x) that is perfectly disentangled w.r.t to z. Then there exists an equivalent generative model with latents $\hat{z} = f(z)$ which is completely entangled and $p(z) = p(\hat{z})$ almost everywhere and thus the same martial distributions P(x). Therefore, the disentanglement method cannot distinguish between the two generative models.

Implications of the Impossibility result

Unsupervised disentanglement is impossible. But if you have

- inductive biases structure can be exploited for better disentanglement.

Visual Proof: Which way is the baseball going? From this evidence alone there are multiple paths the baseball could take and each corresponds to a different generative model. The path of the ball cannot be disentangled without inductiv biases -i.e. flow of time.

What if we had a different definition of Disentanglement?

Definition (Informal Higgins et. Al 2018): A vector representation is called a disentangled representation to a particular decomposition of a symmetry group into subgroups, if it decomposes into independent subspaces, where each subspace is affected by the action of a single subgroup and the action of all other subgroups leave the subspace unaffected.

Groups

- 1. $e \in G$ Identity
- $2.(a \circ b) \circ c = a \circ (b \circ c)$

Associativity

3. $\forall a \in G \quad \exists b \in G$ $a \circ b = e$

Unique Inverses

Symmetries in ML

Shift

Translation Invariance in image labels

Symmetries in ML

Permutation Invariance in Node Labels in a Graph

Symmetries in ML

Rotation

Rotation Equivariance in image features

Example: The dihedral group D_4

e

ľ

 r^2

rs

S

 $r^2 s$

$srs = r^{-1}$

Basic facts on Representations

1. Two representation R and R' are said to be equivalent if $R'(g) = UR(g)U^{\dagger}$ for some Unitary Matrix U

2. A representation R is said to be (completely) reducible if

$$R(g) = U\left(\begin{array}{c} R_1(g) \\ -\end{array}\right)$$

$$\left(R_{2}(g) \right) U^{\dagger}$$

Complete Reducibility

space V. If R fixes the subspace W, then it also fixes W^{\perp} .

the group is finite, and Peter-Weyl (part 2) for continuous.

Theorem. Let R be a representation of a compact group G on a vector

Corollary. Any representation of a compact group is reducible into a direct sum of irreducible representations. This is Maschke's Theorem in

R_2	R_3	R_4
(1)	(1)	$\begin{pmatrix} 1 & 0 \\ 0 & 1 \end{pmatrix}$
(-1)	(-1)	$\begin{pmatrix} i & 0 \\ 0 & -i \end{pmatrix}$
(1)	(1)	$\begin{pmatrix} -1 & 0 \\ 0 & -1 \end{pmatrix}$
(-1)	(-1)	$\begin{pmatrix} -i & 0 \\ 0 & i \end{pmatrix}$
(1)	(-1)	$\begin{pmatrix} 0 & 1 \\ 1 & 0 \end{pmatrix}$
(-1)	(1)	$\begin{pmatrix} 0 & i \\ -i & 0 \end{pmatrix}$
(1)	(-1)	$\begin{pmatrix} 0 & -1 \\ -1 & 0 \end{pmatrix}$
(-1)	(1)	$\begin{pmatrix} 0 & -i \\ i & 0 \end{pmatrix}$

Conjecture

Joey's Principle: Generative Models should focus on learning symmetries in data for disentangled representations.

Free Research Idea: We can now solve the Anomaly Detection problem if we use symmetry based disentangled representations. If a data point is not represented as irreducible representations we simply decompose such that it is. We can then build any density model on these representations.

Theorem 1. (Caselles-Dupré et. Al 2019) (Paraphrased): Symmetry based Disentanglement needs interaction with the environment

The False Dichotomy of Generative Models

Data

Symmetry Group e.g. invariances, equivariances

Geometric Structure e.g. **Riemannian manifold**

Dataset curation: Size, diversity, and temporal

Generative Model

Model Class — i.e. Implicit, Prescribed, Fully Observed

Learning Principle e.g. Variational, MLE, Contrastive

Algorithms e.g. VAE, GANs, RBM

The False Dichotomy of Generative Models

Data

Symmetry Group e.g. invariances, equivariances

Geometric Structure e.g. **Riemannian manifold**

Dataset curation: Size, diversity, and temporal

Generative Model

Model Class — i.e. Implicit, Prescribed, Fully Observed

Learning Principle e.g. Variational, MLE, Contrastive

Algorithms e.g. VAE, GANs, RBM

The False Dichotomy of Generative Models

Data

Symmetry Group e.g. invariances, equivariances

Geometric Structure e.g. **Riemannian manifold**

Dataset curation: Size, diversity, and temporal

Generative Model

Model Class — i.e. Implicit, Prescribed, Fully Observed

Learning Principle e.g. Variational, MLE, Contrastive

Algorithms e.g. VAE, GANs, RBM

Benefits of modelling data in a Generative Model Example with Hyperbolic Geometry:

Embedding Hierarchies in Euclidean space

respect graph distance!

We quickly run out of space! Node Embedding Distance does not

Fig credit: <u>https://openreview.net/pdf?id=BJg73xHtvr</u>

Embedding Hierarchies in Hyperbolic Space

Density Estimation on \mathbb{H}_{K}^{n} Checkerboard

Wrapped Gaussian

2D Spiral

