

Variational Mixture of Normalizing Flows

Guilherme Grijó Pen Freitas Pires

November 13, 2019

Thesis to obtain the Master of Science degree in
Electrical and Computer Engineering



This work in one sentence:

- The development of a mixture of *normalizing flows* and a (variational) training procedure for it.

Table of Contents

- 1 Introduction and Motivation
- 2 Probabilistic Modelling
- 3 Variational Inference
- 4 Normalizing Flows
- 5 Variational Mixture of Normalizing Flows
- 6 Conclusions

Introduction and Motivation

- Deep generative models have been an active area of research, with promising results.

Introduction and Motivation

- Deep generative models have been an active area of research, with promising results.
 - Implicit distributions: Generative adversarial networks [Goodfellow et al., 2014], Variational Autoencoder [Kingma and Welling, 2014]
 - Don't allow explicit access to the density function

Introduction and Motivation

- Deep generative models have been an active area of research, with promising results.
 - Implicit distributions: Generative adversarial networks [Goodfellow et al., 2014], Variational Autoencoder [Kingma and Welling, 2014]
 - Don't allow explicit access to the density function
 - Explicit distributions: Normalizing flows [Rezende and Mohamed, 2015]
 - Allow explicit access to the density function

Introduction and Motivation

- This work

Introduction and Motivation

- This work
 - How to endow normalizing flows with discrete structure?

Introduction and Motivation

- This work
 - How to endow normalizing flows with discrete structure?
 - Or, how to endow mixture models with more expressiveness/flexibility?

- Introductory concepts on probabilistic modelling

Outline

- Introductory concepts on probabilistic modelling
- Variational Inference

Outline

- Introductory concepts on probabilistic modelling
- Variational Inference
 - The chosen method for optimizing the model proposed in this work

Outline

- Introductory concepts on probabilistic modelling
- Variational Inference
 - The chosen method for optimizing the model proposed in this work
- Normalizing Flows

- Introductory concepts on probabilistic modelling
- Variational Inference
 - The chosen method for optimizing the model proposed in this work
- Normalizing Flows
 - The centerpiece of the proposed model

- Introductory concepts on probabilistic modelling
- Variational Inference
 - The chosen method for optimizing the model proposed in this work
- Normalizing Flows
 - The centerpiece of the proposed model
- Variational Mixture of Normalizing Flows

- Introductory concepts on probabilistic modelling
- Variational Inference
 - The chosen method for optimizing the model proposed in this work
- Normalizing Flows
 - The centerpiece of the proposed model
- Variational Mixture of Normalizing Flows
- Experiments and results

- Introductory concepts on probabilistic modelling
- Variational Inference
 - The chosen method for optimizing the model proposed in this work
- Normalizing Flows
 - The centerpiece of the proposed model
- Variational Mixture of Normalizing Flows
- Experiments and results
- Conclusions and future work

Table of Contents

- 1 Introduction and Motivation
- 2 Probabilistic Modelling**
- 3 Variational Inference
- 4 Normalizing Flows
- 5 Variational Mixture of Normalizing Flows
- 6 Conclusions

Probabilistic Modelling: Goal

Given data, find the probability distribution (commonly referred to as the *model*) that is the closest possible approximation to the true distribution of the data.

Probabilistic Modelling: Goal

Informally, via Bayes' Law:

$$p(\text{hypothesis}|\text{data}) = \frac{p(\text{data}|\text{hypothesis})p(\text{hypothesis})}{p(\text{data})}.$$

Probabilistic Modelling: Goal

Informally, via Bayes' Law:

$$p(\text{hypothesis}|\text{data}) = \frac{p(\text{data}|\text{hypothesis})p(\text{hypothesis})}{p(\text{data})}.$$

The goal of probabilistic modelling is to find the optimal hypothesis that maximizes (some form of) this expression.

Probabilistic Modelling: Hypothesis

There are infinite candidate distributions to model the data. In practice, the scope is narrowed to a class of hypothesis.

Probabilistic Modelling: Hypothesis

There are infinite candidate distributions to model the data. In practice, the scope is narrowed to a class of hypothesis.

- Parametric families: $p(\boldsymbol{x}|\boldsymbol{\theta})$

Probabilistic Modelling: Hypothesis

There are infinite candidate distributions to model the data. In practice, the scope is narrowed to a class of hypothesis.

- Parametric families: $p(\mathbf{x}|\boldsymbol{\theta})$
- Particular factorizations: e.g. $p(\mathbf{x}) = \prod_i^N p(x_i|x_{i-1})$

Probabilistic Modelling: Hypothesis

There are infinite candidate distributions to model the data. In practice, the scope is narrowed to a class of hypothesis.

- Parametric families: $p(\mathbf{x}|\boldsymbol{\theta})$
- Particular factorizations: e.g. $p(\mathbf{x}) = \prod_i^N p(x_i|x_{i-1})$
- Latent variables: $p(\mathbf{x}) = \int p(\mathbf{x}, \mathbf{z})d\mathbf{z}$ ¹

¹In the case of discrete latent variables, the integral becomes a sum

Probabilistic Modelling: Hypothesis

There are infinite candidate distributions to model the data. In practice, the scope is narrowed to a class of hypothesis.

- Parametric families: $p(\mathbf{x}|\boldsymbol{\theta})$
- Particular factorizations: e.g. $p(\mathbf{x}) = \prod_i^N p(x_i|x_{i-1})$
- Latent variables: $p(\mathbf{x}) = \int p(\mathbf{x}, \mathbf{z})d\mathbf{z}$ ¹
- All of the above, combined

¹In the case of discrete latent variables, the integral becomes a sum

Probabilistic Modelling: Hypothesis

Example: Mixture Models.

Probabilistic Modelling: Hypothesis

Example: Mixture Models.

A mixture model is defined as:

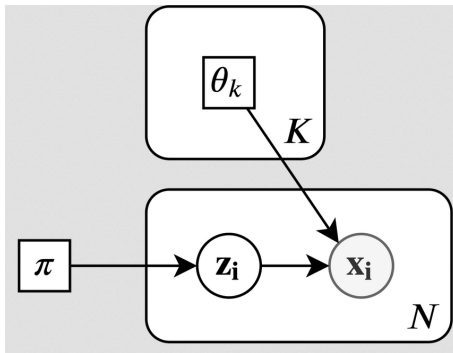
$$p(\mathbf{x}, \mathbf{z}) = p(\mathbf{x}|\mathbf{z}, \boldsymbol{\theta})p(\mathbf{z})$$

Probabilistic Modelling: Hypothesis

Example: Mixture Models.

A mixture model is defined as:

$$p(\mathbf{x}, \mathbf{z}) = p(\mathbf{x}|\mathbf{z}, \boldsymbol{\theta})p(\mathbf{z})$$



Probabilistic Modelling: Parameter Estimation

Given data x and a parametric model $p(x|\theta)$, to estimate θ :

Probabilistic Modelling: Parameter Estimation

Given data x and a parametric model $p(x|\theta)$, to estimate θ :

Maximum-likelihood:

$$\begin{cases} \hat{\theta}_{ML} = \operatorname{argmax}_{\theta} \mathcal{L}(\theta), \\ \mathcal{L}(\theta) = p(x|\theta) \end{cases}$$

Probabilistic Modelling: Parameter Estimation

Given data \mathbf{x} and a parametric model $p(\mathbf{x}|\boldsymbol{\theta})$, to estimate $\boldsymbol{\theta}$:

Maximum-likelihood:

$$\begin{cases} \hat{\boldsymbol{\theta}}_{ML} = \operatorname{argmax}_{\boldsymbol{\theta}} \mathcal{L}(\boldsymbol{\theta}), \\ \mathcal{L}(\boldsymbol{\theta}) = p(\mathbf{x}|\boldsymbol{\theta}) \end{cases}$$

Maximum a posteriori:

$$\begin{cases} \hat{\boldsymbol{\theta}}_{MAP} = \operatorname{argmax}_{\boldsymbol{\theta}} p(\boldsymbol{\theta}|\mathbf{x}), \\ p(\boldsymbol{\theta}|\mathbf{x}) = \frac{p(\mathbf{x}|\boldsymbol{\theta})p(\boldsymbol{\theta})}{p(\mathbf{x})}. \end{cases}$$

Probabilistic Modelling: Inference

Given data, x , and a model $p(x, z)$, one is generally interested in finding the posterior:

Probabilistic Modelling: Inference

Given data, \mathbf{x} , and a model $p(\mathbf{x}, \mathbf{z})$, one is generally interested in finding the posterior:

$$\begin{aligned} p(\mathbf{z}|\mathbf{x}) &= \frac{p(\mathbf{x}|\mathbf{z})p(\mathbf{z})}{p(\mathbf{x})} \\ &= \frac{p(\mathbf{x}|\mathbf{z})p(\mathbf{z})}{\int p(\mathbf{x}|\mathbf{z}')p(\mathbf{z}')d\mathbf{z}'} \end{aligned}$$

Probabilistic Modelling: Inference

Given data, \mathbf{x} , and a model $p(\mathbf{x}, \mathbf{z})$, one is generally interested in finding the posterior:

$$\begin{aligned} p(\mathbf{z}|\mathbf{x}) &= \frac{p(\mathbf{x}|\mathbf{z})p(\mathbf{z})}{p(\mathbf{x})} \\ &= \frac{p(\mathbf{x}|\mathbf{z})p(\mathbf{z})}{\int p(\mathbf{x}|\mathbf{z}')p(\mathbf{z}')d\mathbf{z}'} \end{aligned}$$

In general, the integral in the denominator is intractable.

Probabilistic Modelling: Inference

Given data, \mathbf{x} , and a model $p(\mathbf{x}, \mathbf{z})$, one is generally interested in finding the posterior:

$$\begin{aligned} p(\mathbf{z}|\mathbf{x}) &= \frac{p(\mathbf{x}|\mathbf{z})p(\mathbf{z})}{p(\mathbf{x})} \\ &= \frac{p(\mathbf{x}|\mathbf{z})p(\mathbf{z})}{\int p(\mathbf{x}|\mathbf{z}')p(\mathbf{z}')d\mathbf{z}'} \end{aligned}$$

In general, the integral in the denominator is intractable.

→ Approximate inference is required.

Table of Contents

- 1 Introduction and Motivation
- 2 Probabilistic Modelling
- 3 Variational Inference**
- 4 Normalizing Flows
- 5 Variational Mixture of Normalizing Flows
- 6 Conclusions

Variational Inference: Goal

Variational Inference is one way of dealing with the intractability previously described.

Variational Inference: Goal

Variational Inference is one way of dealing with the intractability previously described.

Given a *variational* family $q(\mathbf{z}; \boldsymbol{\lambda})$, find the parameters $\boldsymbol{\lambda}$ that minimize the Kullback-Leibler divergence between $q(\mathbf{z}; \boldsymbol{\lambda})$ and $p(\mathbf{z}|\mathbf{x})$

Variational Inference: ELBO

$$KL(q||p) = \int q(\mathbf{z}) \log \frac{q(\mathbf{z})}{p(\mathbf{z})} d\mathbf{z}$$

Variational Inference: ELBO

$$\begin{aligned}KL(q||p) &= \int q(\mathbf{z}) \log \frac{q(\mathbf{z})}{p(\mathbf{z})} d\mathbf{z} \\ &= \int q(\mathbf{z}) (\log q(\mathbf{z}) - \log p(\mathbf{z}|\mathbf{x})) d\mathbf{z}\end{aligned}$$

Variational Inference: ELBO

$$\begin{aligned}KL(q||p) &= \int q(\mathbf{z}) \log \frac{q(\mathbf{z})}{p(\mathbf{z})} d\mathbf{z} \\&= \int q(\mathbf{z}) (\log q(\mathbf{z}) - \log p(\mathbf{z}|\mathbf{x})) d\mathbf{z} \\&= \int q(\mathbf{z}) (\log q(\mathbf{z}) - (\log p(\mathbf{x}, \mathbf{z}) - \log p(\mathbf{x}))) d\mathbf{z}\end{aligned}$$

Variational Inference: ELBO

$$\begin{aligned}KL(q||p) &= \int q(\mathbf{z}) \log \frac{q(\mathbf{z})}{p(\mathbf{z})} d\mathbf{z} \\&= \int q(\mathbf{z}) (\log q(\mathbf{z}) - \log p(\mathbf{z}|\mathbf{x})) d\mathbf{z} \\&= \int q(\mathbf{z}) (\log q(\mathbf{z}) - (\log p(\mathbf{x}, \mathbf{z}) - \log p(\mathbf{x}))) d\mathbf{z} \\&= \mathbb{E}_q[\log q(\mathbf{z})] - \mathbb{E}_q[\log p(\mathbf{x}, \mathbf{z})] + \log p(\mathbf{x})\end{aligned}$$

Variational Inference: ELBO

$$\begin{aligned}KL(q||p) &= \int q(\mathbf{z}) \log \frac{q(\mathbf{z})}{p(\mathbf{z})} d\mathbf{z} \\&= \int q(\mathbf{z}) (\log q(\mathbf{z}) - \log p(\mathbf{z}|\mathbf{x})) d\mathbf{z} \\&= \int q(\mathbf{z}) (\log q(\mathbf{z}) - (\log p(\mathbf{x}, \mathbf{z}) - \log p(\mathbf{x}))) d\mathbf{z} \\&= \mathbb{E}_q[\log q(\mathbf{z})] - \mathbb{E}_q[\log p(\mathbf{x}, \mathbf{z})] + \log p(\mathbf{x})\end{aligned}$$

Which yields the lower bound (ELBO):

$$\begin{aligned}ELBO(q) &= \mathbb{E}_q[\log p(\mathbf{x}, \mathbf{z})] - \mathbb{E}_q[\log q(\mathbf{z})] \\&= \mathbb{E}_q[\log p(\mathbf{x}|\mathbf{z})] + \mathbb{E}_q[\log p(\mathbf{z})] - \mathbb{E}_q[\log q(\mathbf{z})]\end{aligned}$$

Table of Contents

- 1 Introduction and Motivation
- 2 Probabilistic Modelling
- 3 Variational Inference
- 4 Normalizing Flows**
- 5 Variational Mixture of Normalizing Flows
- 6 Conclusions

Normalizing Flows: Change of Variables

Given

$$\begin{cases} z \sim p(z) \\ x = g(z; \theta) \end{cases}$$

Normalizing Flows: Change of Variables

Given

$$\begin{cases} z \sim p(z) \\ x = g(z; \theta) \end{cases}$$

then:

$$\begin{aligned} f_X(\mathbf{x}) &= f_Z(g^{-1}(\mathbf{x}; \theta)) \left| \det \left(\frac{d}{d\mathbf{x}} g^{-1}(\mathbf{x}; \theta) \right) \right| \\ &= f_Z(g^{-1}(\mathbf{x}; \theta)) \left| \det \left(\frac{d}{d\mathbf{z}} g(\mathbf{z}; \theta) \right) \Big|_{\mathbf{z}=g^{-1}(\mathbf{x}; \theta)} \right|^{-1} \end{aligned}$$

Normalizing Flows: Change of Variables

Given

$$\begin{cases} z \sim p(z) \\ x = g(z; \theta) \end{cases}$$

then:

$$\begin{aligned} f_X(x) &= f_Z(g^{-1}(x; \theta)) \left| \det \left(\frac{d}{dx} g^{-1}(x; \theta) \right) \right| \\ &= f_Z(g^{-1}(x; \theta)) \left| \det \left(\frac{d}{dz} g(z; \theta) \right) \Big|_{z=g^{-1}(x; \theta)} \right|^{-1} \end{aligned}$$

This can be optimize w.r.t. θ so as to approximate an arbitrary distribution

Normalizing Flows: Change of Variables

The above can be useful if

Normalizing Flows: Change of Variables

The above can be useful if

- The base density has a closed form expression and is easy to sample from

Normalizing Flows: Change of Variables

The above can be useful if

- The base density has a closed form expression and is easy to sample from
- The determinant of the Jacobian of g is computationally cheap - not the case, in general

Normalizing Flows: Change of Variables

The above can be useful if

- The base density has a closed form expression and is easy to sample from
- The determinant of the Jacobian of g is computationally cheap - not the case, in general
- The gradient of the determinant of the Jacobian of g w.r.t θ is computationally cheap

Normalizing Flows: Change of Variables

The framework of Normalizing Flows consists of composing several transformations that fulfill the three listed conditions.

Normalizing Flows: Affine Coupling Layer

An example of such a transformation is the Affine Coupling Layer [Dinh, Sohl-Dickstein, and Bengio, 2017].

Splitting z into (z_1, z_2) ,

Normalizing Flows: Affine Coupling Layer

An example of such a transformation is the Affine Coupling Layer [Dinh, Sohl-Dickstein, and Bengio, 2017].

Splitting z into (z_1, z_2) ,

$$\begin{cases} x_1 &= z_1 \odot \exp(s(z_2)) + t(z_2) \\ x_2 &= z_2. \end{cases}$$

Normalizing Flows: Affine Coupling Layer

An example of such a transformation is the Affine Coupling Layer [Dinh, Sohl-Dickstein, and Bengio, 2017].

Splitting z into (z_1, z_2) ,

$$\begin{cases} x_1 &= z_1 \odot \exp(s(z_2)) + t(z_2) \\ x_2 &= z_2. \end{cases}$$

This transformation has the following Jacobian matrix:

$$J_{f(z)} = \begin{bmatrix} \frac{\partial x_1}{\partial z_1} & \frac{\partial x_1}{\partial z_2} \\ \frac{\partial x_2}{\partial z_1} & \frac{\partial x_2}{\partial z_2} \end{bmatrix} = \begin{bmatrix} \text{diag}(\exp(s(z_2))) & \frac{\partial x_1}{\partial z_2} \\ \mathbf{0} & I \end{bmatrix}$$

Table of Contents

- 1 Introduction and Motivation
- 2 Probabilistic Modelling
- 3 Variational Inference
- 4 Normalizing Flows
- 5 Variational Mixture of Normalizing Flows**
- 6 Conclusions

How can we leverage the flexibility of normalizing flows, and endow it with multimodal, discrete structure, like in a mixture model?

How can we leverage the flexibility of normalizing flows, and endow it with multimodal, discrete structure, like in a mixture model?

Mixture of normalizing flows.

How can we leverage the flexibility of normalizing flows, and endow it with multimodal, discrete structure, like in a mixture model?

Mixture of normalizing flows. \rightarrow Approximate inference is required.

Recall the ELBO:

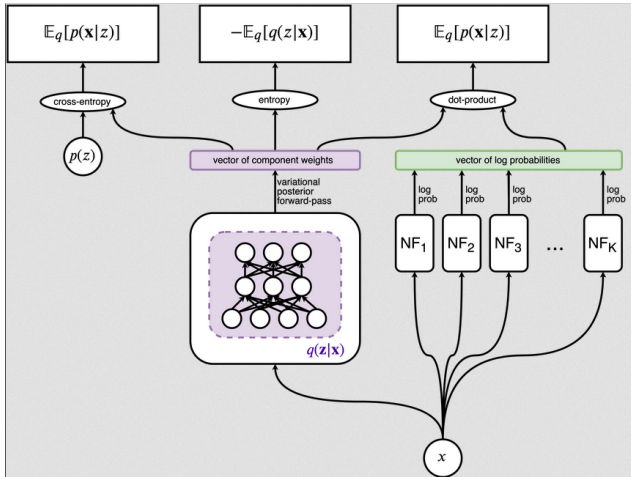
$$ELBO(q) = \mathbb{E}_q[\log p(\boldsymbol{x}|\boldsymbol{z})] + \mathbb{E}_q[\log p(\boldsymbol{z})] - \mathbb{E}_q[\log q(\boldsymbol{z})]$$

Recall the ELBO:

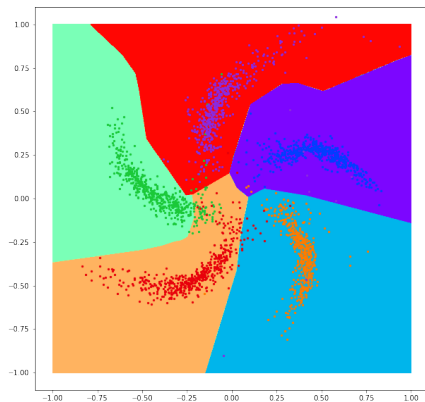
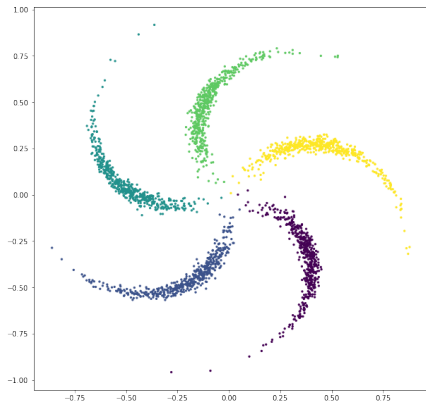
$$ELBO(q) = \mathbb{E}_q[\log p(\mathbf{x}|\mathbf{z})] + \mathbb{E}_q[\log p(\mathbf{z})] - \mathbb{E}_q[\log q(\mathbf{z})]$$

Let the variational posterior $q(\mathbf{z}|\mathbf{x})$ be parameterized by a neural network. We jointly optimize this objective, hence we learn the variational posterior and the generative components simultaneously.

VMoNF: Overview

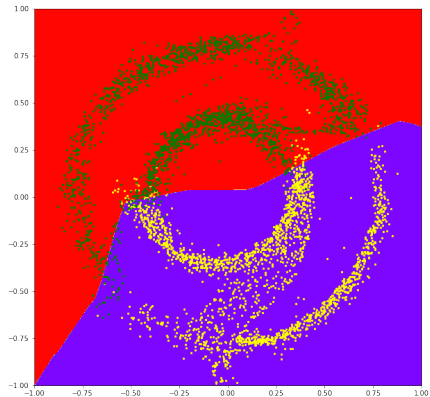
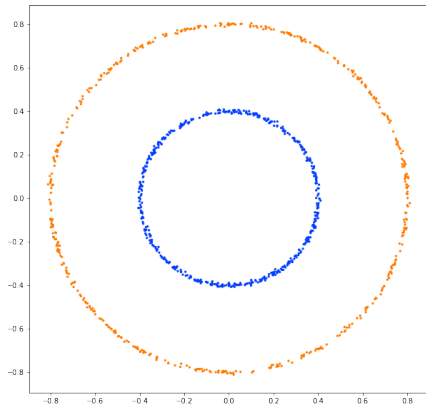


VMoNF: Experiments - Pinwheel (5 wings)

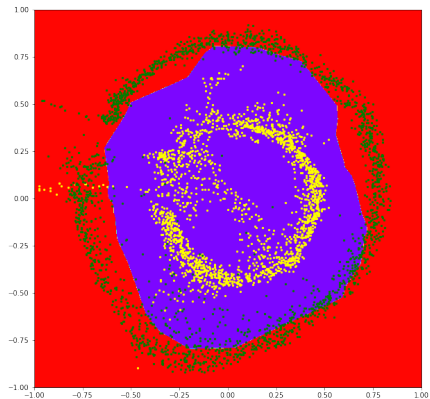
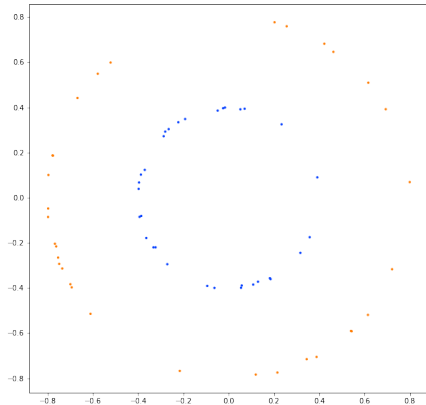


Training Animation

VMoNF: Experiments - 2 Circles



VMoNF: Experiments - 2 Circles (semi supervised)



Note: 32 labeled points, 1024 unlabeled points

VMoNF: Experiments - MNIST

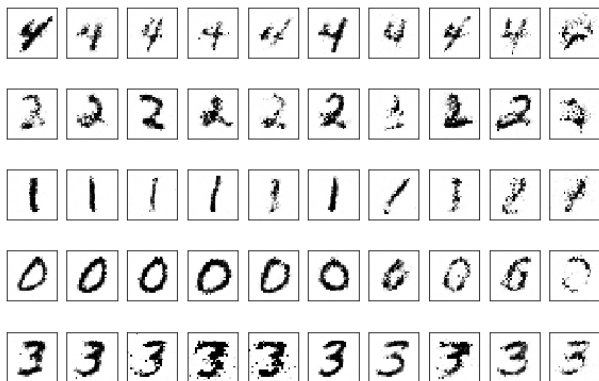


Table of Contents

- 1 Introduction and Motivation
- 2 Probabilistic Modelling
- 3 Variational Inference
- 4 Normalizing Flows
- 5 Variational Mixture of Normalizing Flows
- 6 Conclusions**

Conclusions and Future Work

- Similar work is being pursued and published in prominent venues: [Dinh, Sohl-Dickstein, et al., 2019; Izmailov et al., 2019]

Conclusions and Future Work

- Similar work is being pursued and published in prominent venues: [Dinh, Sohl-Dickstein, et al., 2019; Izmailov et al., 2019]
- Formally describe the reasons why the model fails in cases like the 2 circles \rightarrow Topology

Conclusions and Future Work

- Similar work is being pursued and published in prominent venues: [Dinh, Sohl-Dickstein, et al., 2019; Izmailov et al., 2019]
- Formally describe the reasons why the model fails in cases like the 2 circles → Topology
 - Investigate the effect of a consistency loss regularization term

Conclusions and Future Work

- Similar work is being pursued and published in prominent venues: [Dinh, Sohl-Dickstein, et al., 2019; Izmailov et al., 2019]
- Formally describe the reasons why the model fails in cases like the 2 circles → Topology
 - Investigate the effect of a consistency loss regularization term
- Weight-sharing between components

Conclusions and Future Work

- Similar work is being pursued and published in prominent venues: [Dinh, Sohl-Dickstein, et al., 2019; Izmailov et al., 2019]
- Formally describe the reasons why the model fails in cases like the 2 circles → Topology
 - Investigate the effect of a consistency loss regularization term
- Weight-sharing between components
- Balance between complexities

Conclusions and Future Work

- Similar work is being pursued and published in prominent venues: [Dinh, Sohl-Dickstein, et al., 2019; Izmailov et al., 2019]
- Formally describe the reasons why the model fails in cases like the 2 circles → Topology
 - Investigate the effect of a consistency loss regularization term
- Weight-sharing between components
- Balance between complexities
- (Controlled) component annihilation

Thank you!