

CSCI 699 - ProbGen

Probabilistic and Generative Models

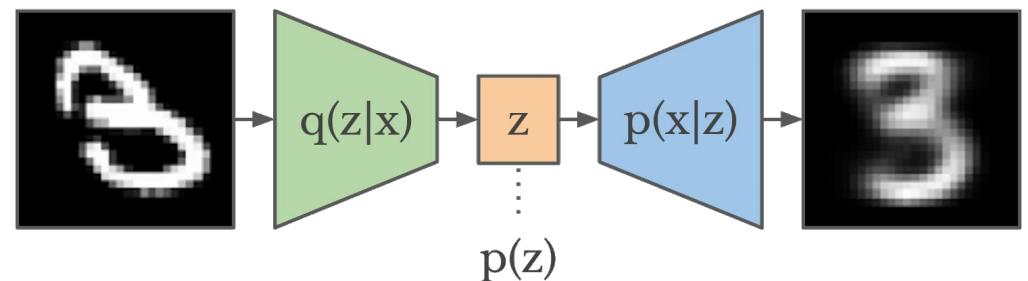
Willie Neiswanger

Lecture 7 - VAEs (continued) and Intro to Score-based Generative Models

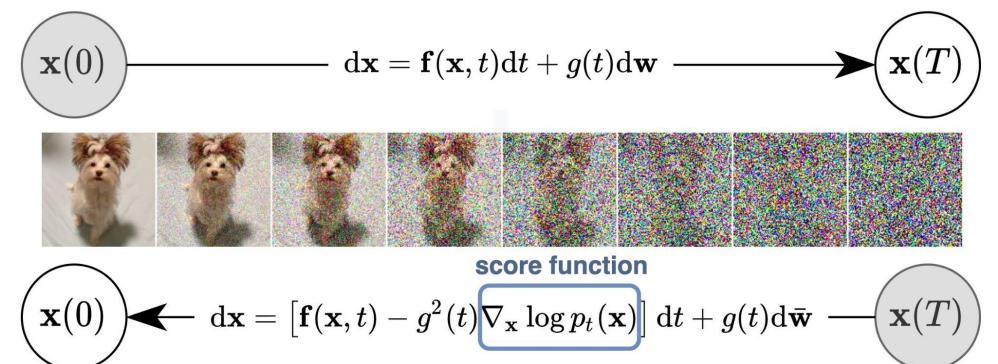
Today

Today

Lecture: Continued explanation of the variational autoencoder (VAE), and then intro to score-based generative models.



"Building Variational Auto-Encoders in TensorFlow", Danijar Hafner

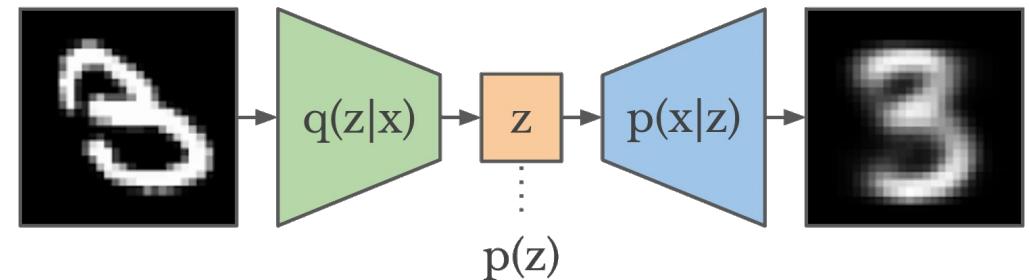


Source: Yang Song, "Generative Modeling by Estimating Gradients of the Data Distribution"

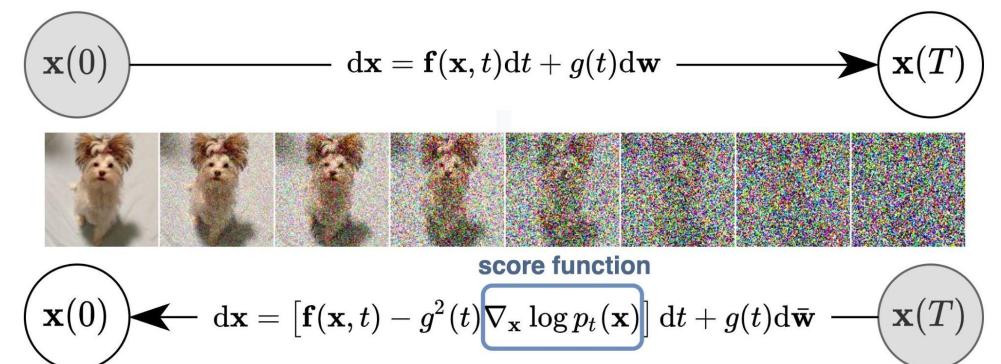
Today

Lecture: Continued explanation of the variational autoencoder (VAE), and then intro to score-based generative models.

- Review: Variational inference (VI) methods
- Variational autoencoders (VAEs)
 - Reparameterization trick, SGVB.
 - Relation to (classic) autoencoders.



"Building Variational Auto-Encoders in TensorFlow", Danijar Hafner

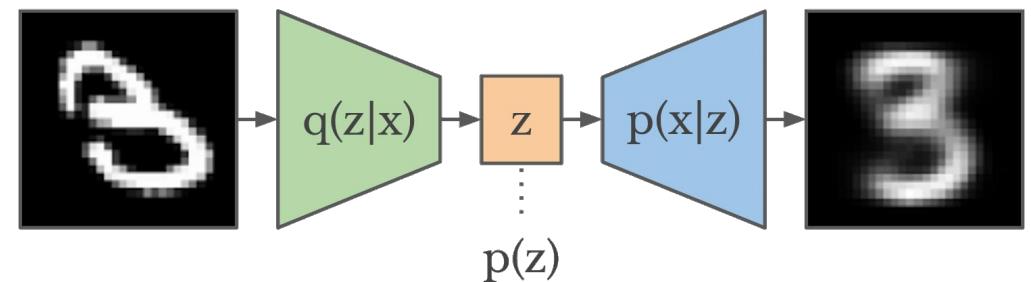


Source: Yang Song, "Generative Modeling by Estimating Gradients of the Data Distribution"

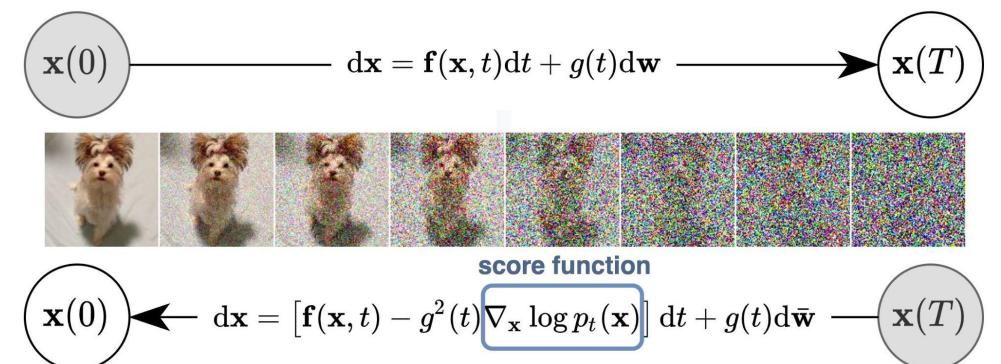
Today

Lecture: Continued explanation of the variational autoencoder (VAE), and then intro to score-based generative models.

- Review: Variational inference (VI) methods
- Variational autoencoders (VAEs)
 - Reparameterization trick, SGVB.
 - Relation to (classic) autoencoders.
- Intro to score-based generative models



"Building Variational Auto-Encoders in TensorFlow", Danijar Hafner

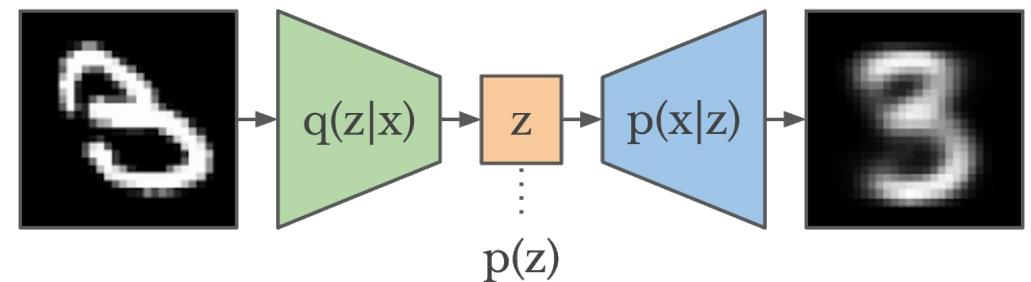


Source: Yang Song, "Generative Modeling by Estimating Gradients of the Data Distribution"

Today

Lecture: Continued explanation of the variational autoencoder (VAE), and then intro to score-based generative models.

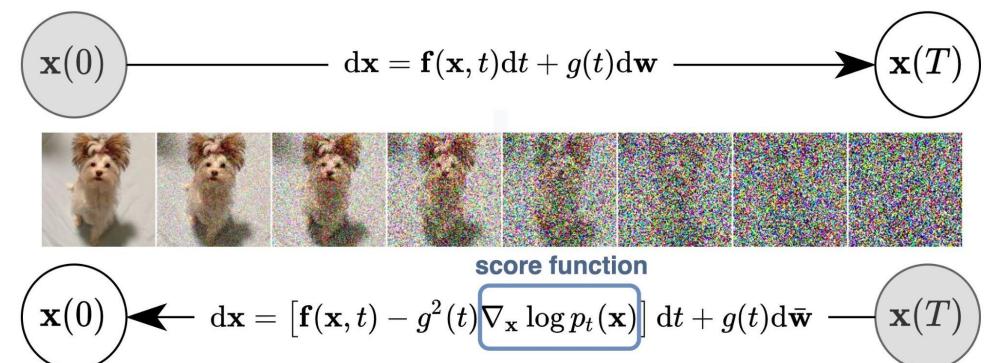
- Review: Variational inference (VI) methods
- Variational autoencoders (VAEs)
 - Reparameterization trick, SGVB.
 - Relation to (classic) autoencoders.
- Intro to score-based generative models



"Building Variational Auto-Encoders in TensorFlow", Danijar Hafner

After:

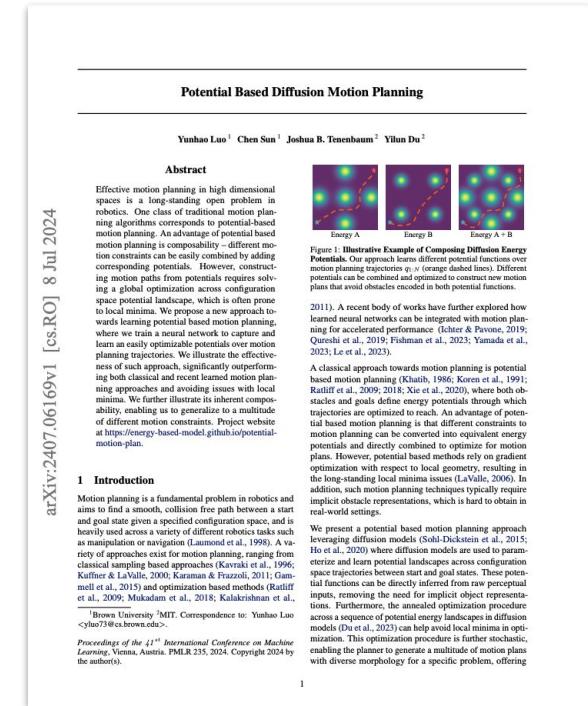
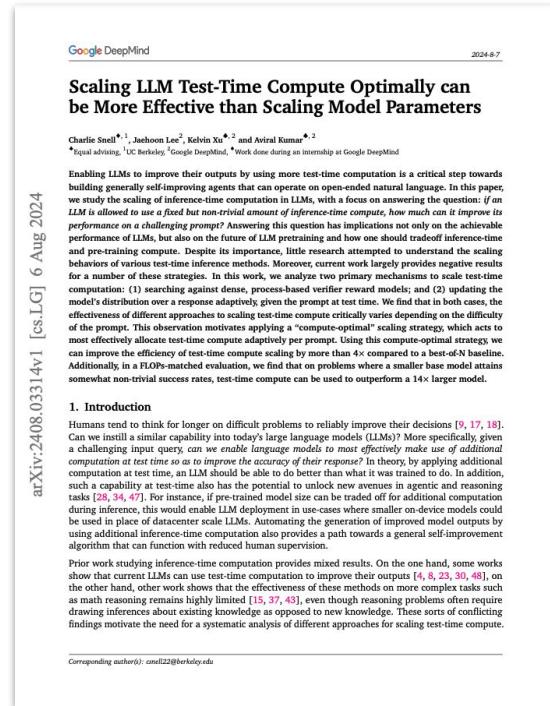
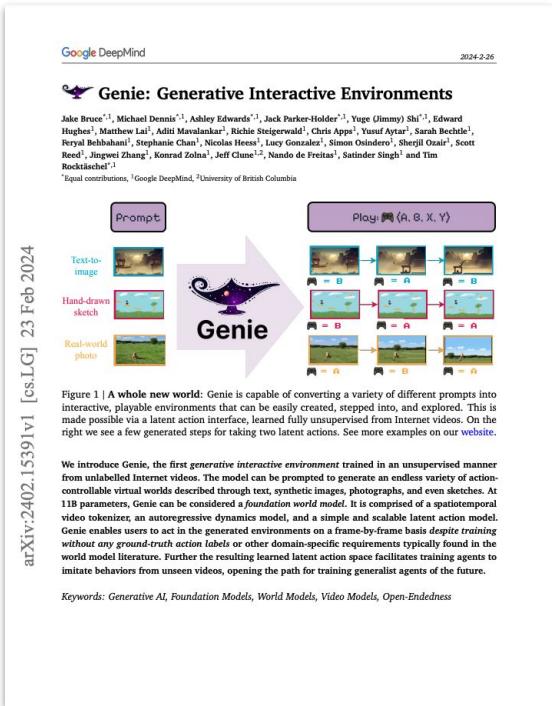
- First groups of four paper presentations.



Source: Yang Song, "Generative Modeling by Estimating Gradients of the Data Distribution"

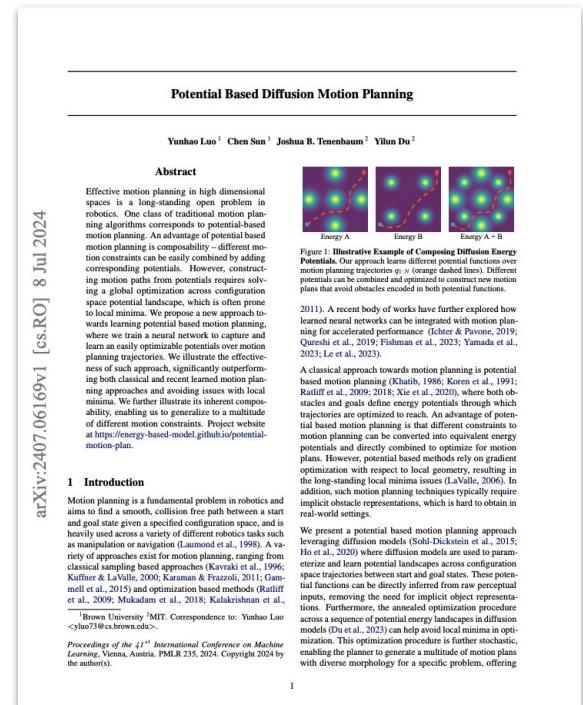
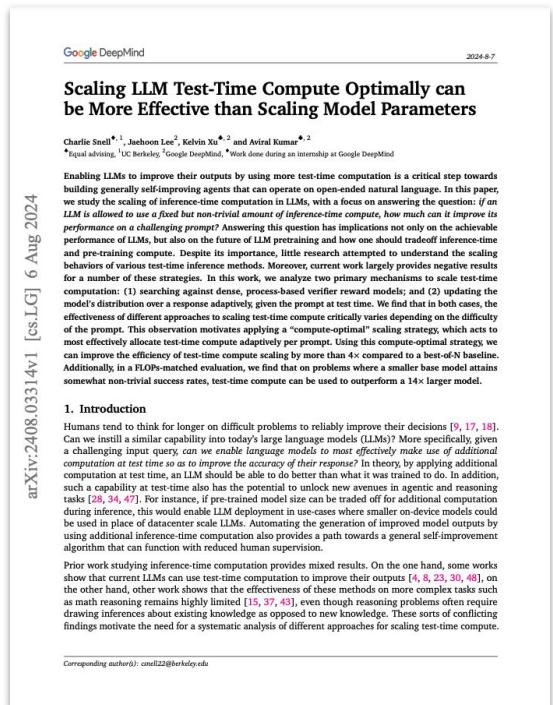
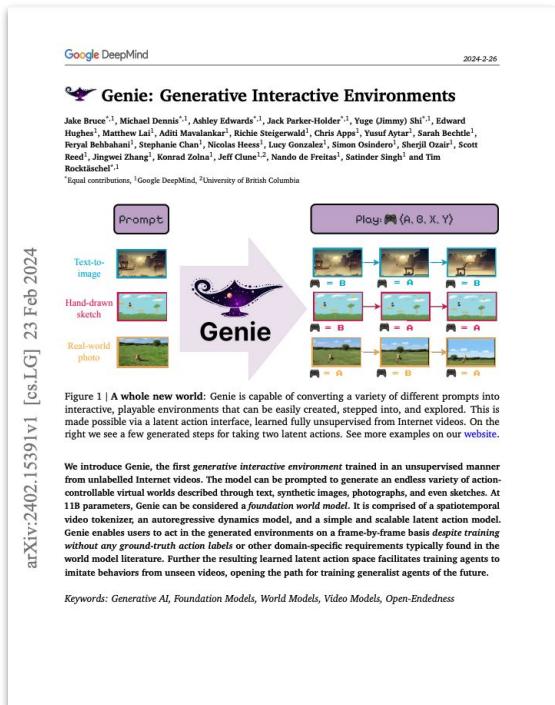
Today

After: First groups of four paper presentations.



Today

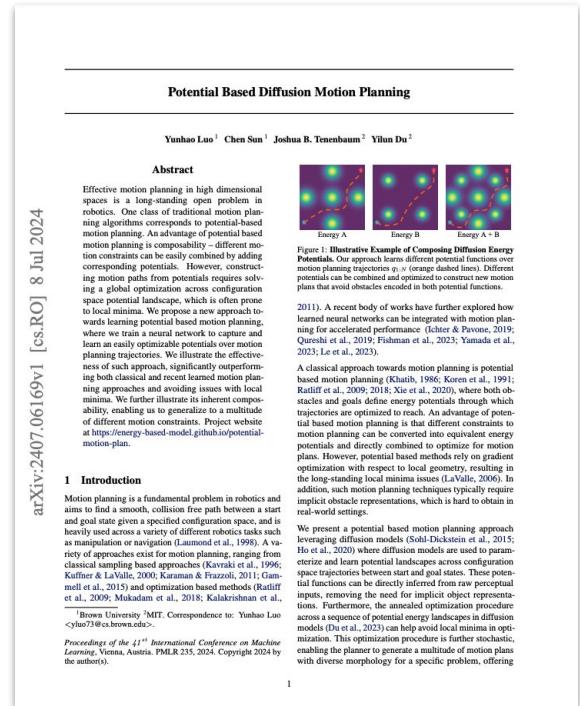
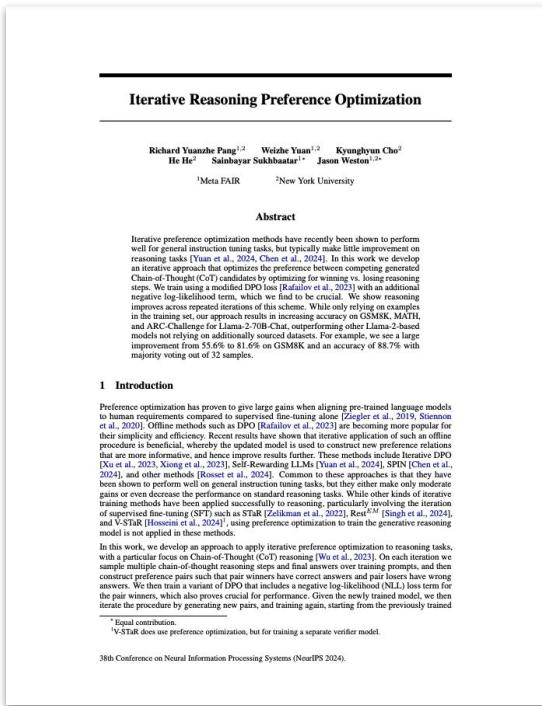
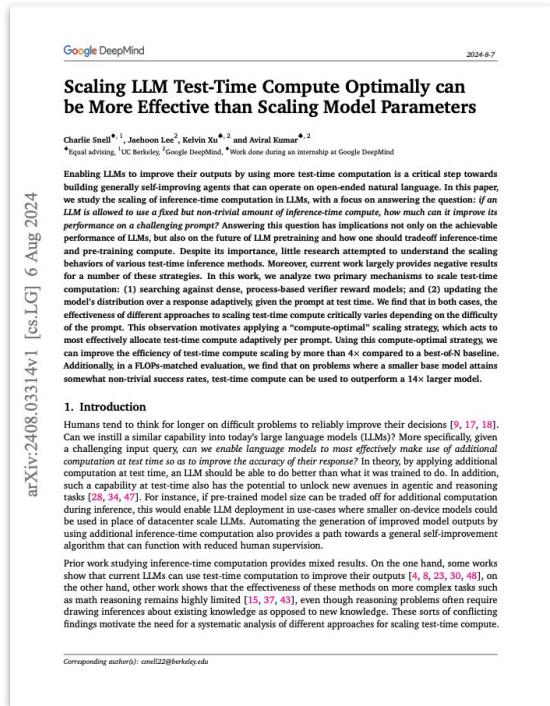
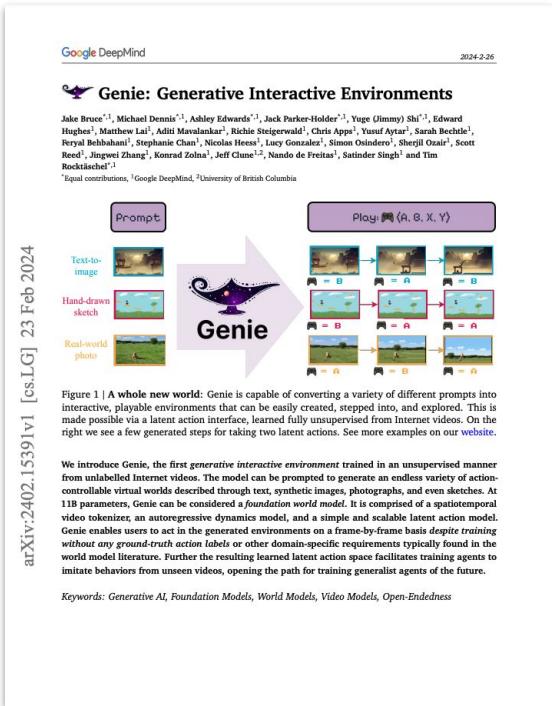
After: First groups of four paper presentations.



Generative Interactive Environments

Today

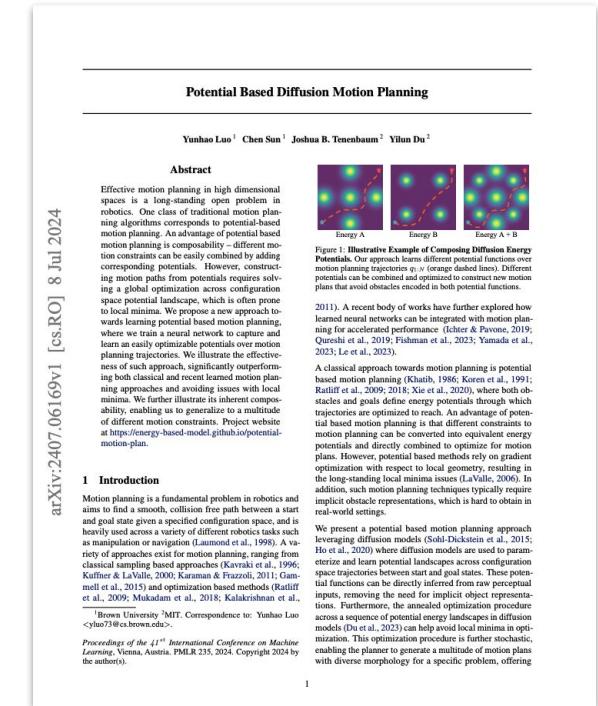
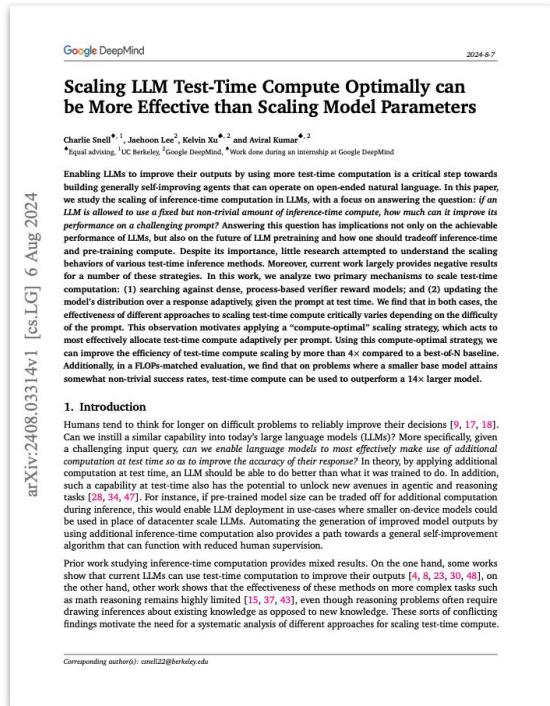
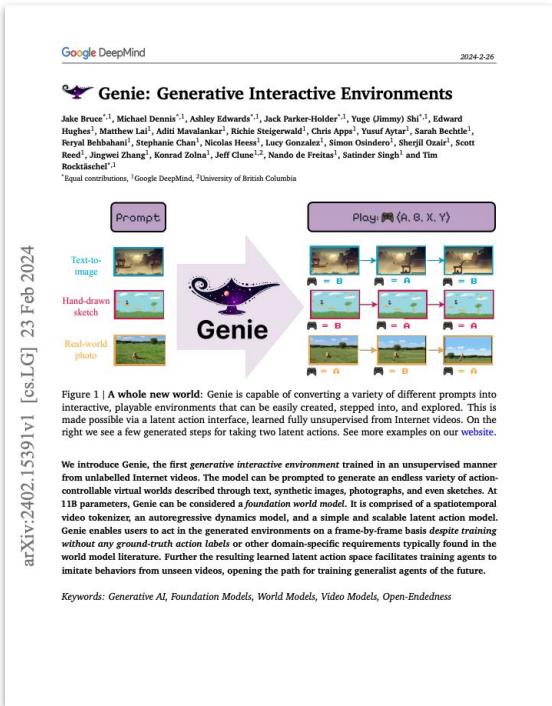
After: First groups of four paper presentations.



Scaling LLM Test-Time Compute

Today

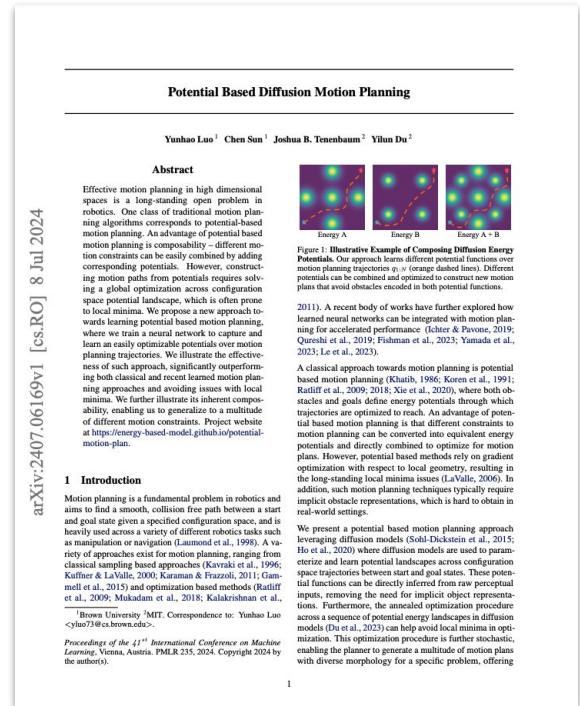
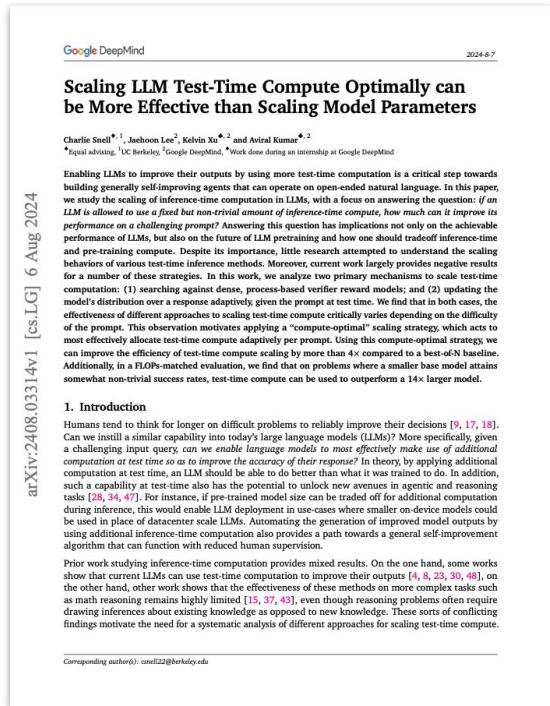
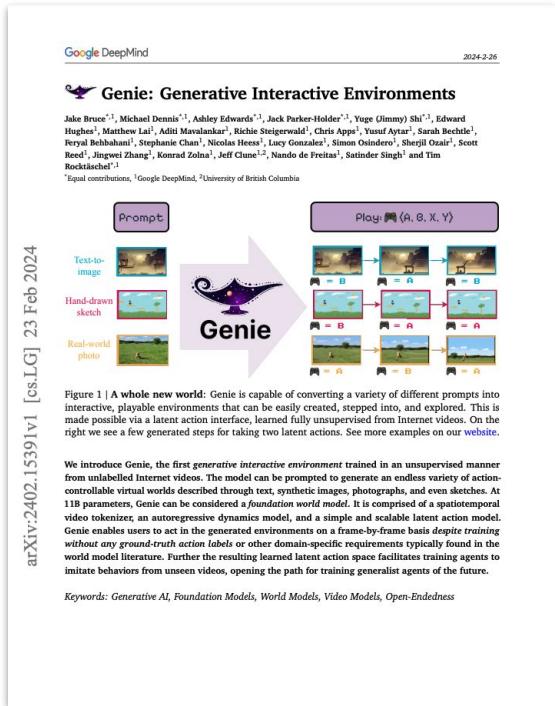
After: First groups of four paper presentations.



Iterative Preference Optimization

Today

After: First groups of four paper presentations.



Diffusion Motion Planning

Paper Presentations

Paper Presentations – Goal

- During the semester, each student gives **one 20-minute presentation** on a paper relevant to probabilistic and generative models – [spreadsheet here](#).
- Ideally from a modern machine learning conference (e.g., NeurIPS, ICML, ICLR, AAAI, etc)
 - Could be a “classic” (older) paper relevant to modern models, as well.
- This will accomplish a few things:
 - Gives the class a chance to see a broad set of interesting probabilistic/generative modeling papers.
 - Gives each student (more) experience with distilling key content from a paper & presenting it.

Paper Presentations – Content of each paper presentation?

Ideally I'd like you to include:

Paper Presentations – Content of each paper presentation?

Ideally I'd like you to include:

1. Motivation (main problem the paper is trying to solve).
2. Brief background on prior work (*i.e.*, related methods, relevant papers).
3. Method and key equations.
 - It'll be important to boil the method down to a concise explanation given ~20 minute presentation.
4. Empirical results – how well does the method work in practice?
5. Your opinion on the paper, its (potential) impact/promise, etc.

Review: Variational Inference

Variational Inference – Overview

Variational Inference – Overview

We can take the strategy/philosophy of *inference as optimization*.

This method of approximate inference is known as **variational inference (VI)**.

Variational Inference – Overview

We can take the strategy/philosophy of *inference as optimization*.

This method of approximate inference is known as **variational inference (VI)**.

The main idea:

Variational Inference – Overview

We can take the strategy/philosophy of *inference as optimization*.

This method of approximate inference is known as **variational inference (VI)**.

The main idea:

- Suppose we are given an intractable (*complex*) probability distribution $p(x)$.

Variational Inference – Overview

We can take the strategy/philosophy of *inference as optimization*.

This method of approximate inference is known as **variational inference (VI)**.

The main idea:

E.g., only known up to a normalization constant

- Suppose we are given an intractable (*complex*) probability distribution $p(x)$.

Variational Inference – Overview

We can take the strategy/philosophy of *inference as optimization*.

This method of approximate inference is known as **variational inference (VI)**.

The main idea:

- Suppose we are given an intractable (*complex*) probability distribution $p(x)$.
- VI aims to solve an optimization problem over a set of tractable distributions Q .

Variational Inference – Overview

We can take the strategy/philosophy of *inference as optimization*.

This method of approximate inference is known as **variational inference (VI)**.

The main idea:

- Suppose we are given an intractable (*complex*) probability distribution $p(x)$.
- VI aims to solve an optimization problem over a set of tractable distributions Q .
- And, after optimization, return a $q(x) \in Q$ that is *most similar* to $p(x)$.

Variational Inference – Overview

We can take the strategy/philosophy of *inference as optimization*.

This method of approximate inference is known as **variational inference (VI)**.

The main idea:

- Suppose we are given an intractable (*complex*) probability distribution $p(x)$.
- VI aims to solve an optimization problem over a set of tractable distributions Q .
- And, after optimization, return a $q(x) \in Q$ that is *most similar* to $p(x)$.

Typically: in terms
of KL divergence

Variational Inference – Overview

We can take the strategy/philosophy of *inference as optimization*.

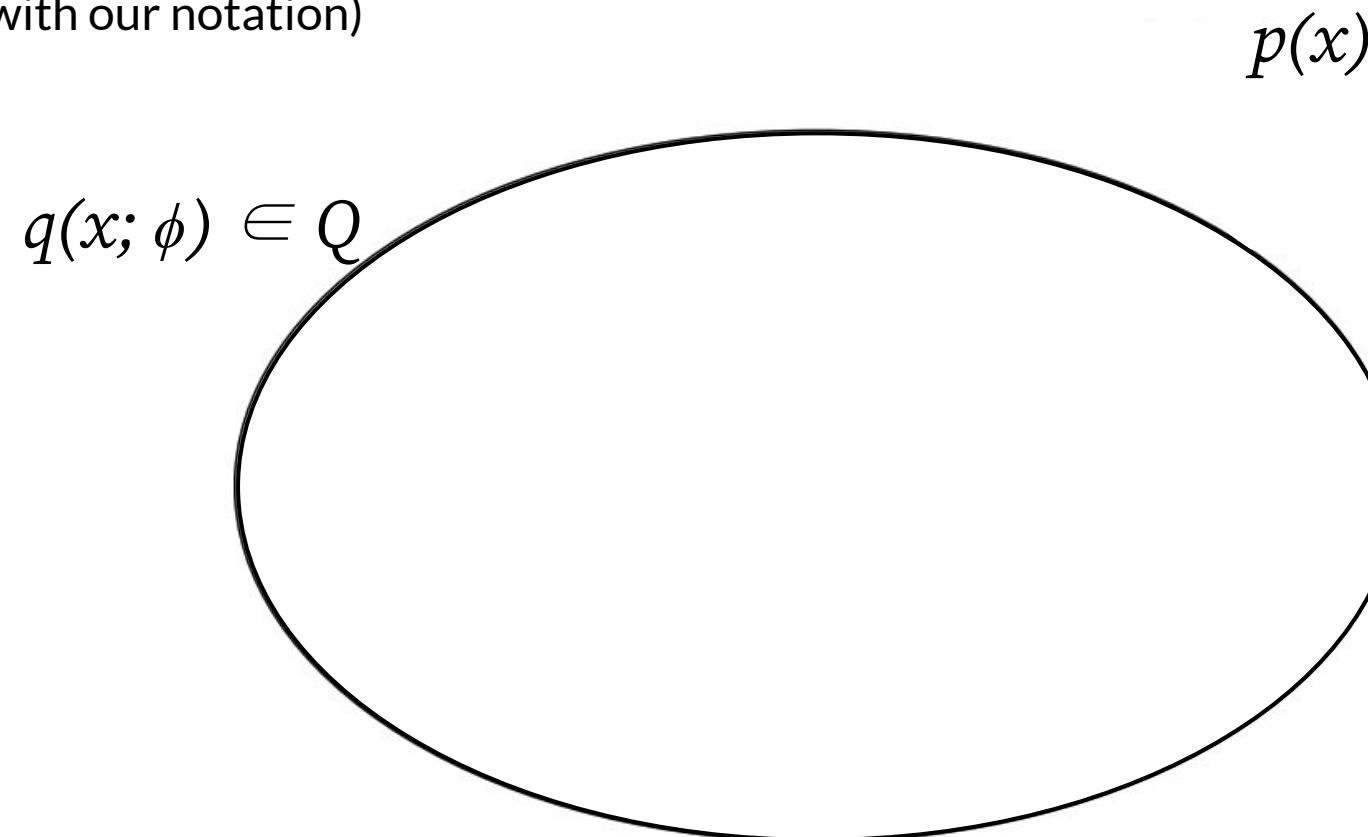
This method of approximate inference is known as **variational inference (VI)**.

The main idea:

- Suppose we are given an intractable (*complex*) probability distribution $p(x)$.
- VI aims to solve an optimization problem over a set of tractable distributions Q .
- And, after optimization, return a $q(x) \in Q$ that is *most similar* to $p(x)$.
- We can then make an inference query on $q(x)$ instead of $p(x)$.

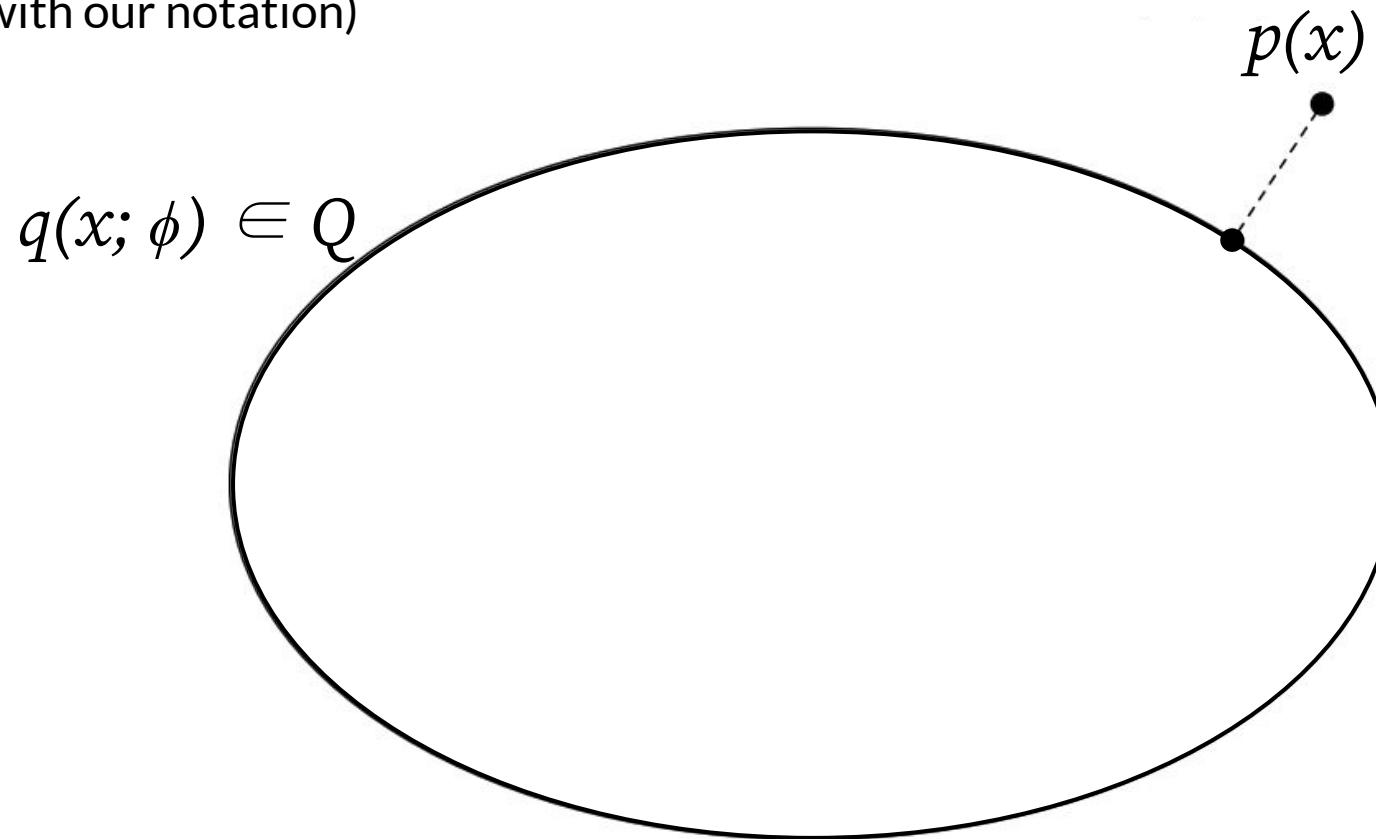
Variational Inference – Overview

An illustration from Dave Blei's Variational Inference Tutorial:
(replaced with our notation)



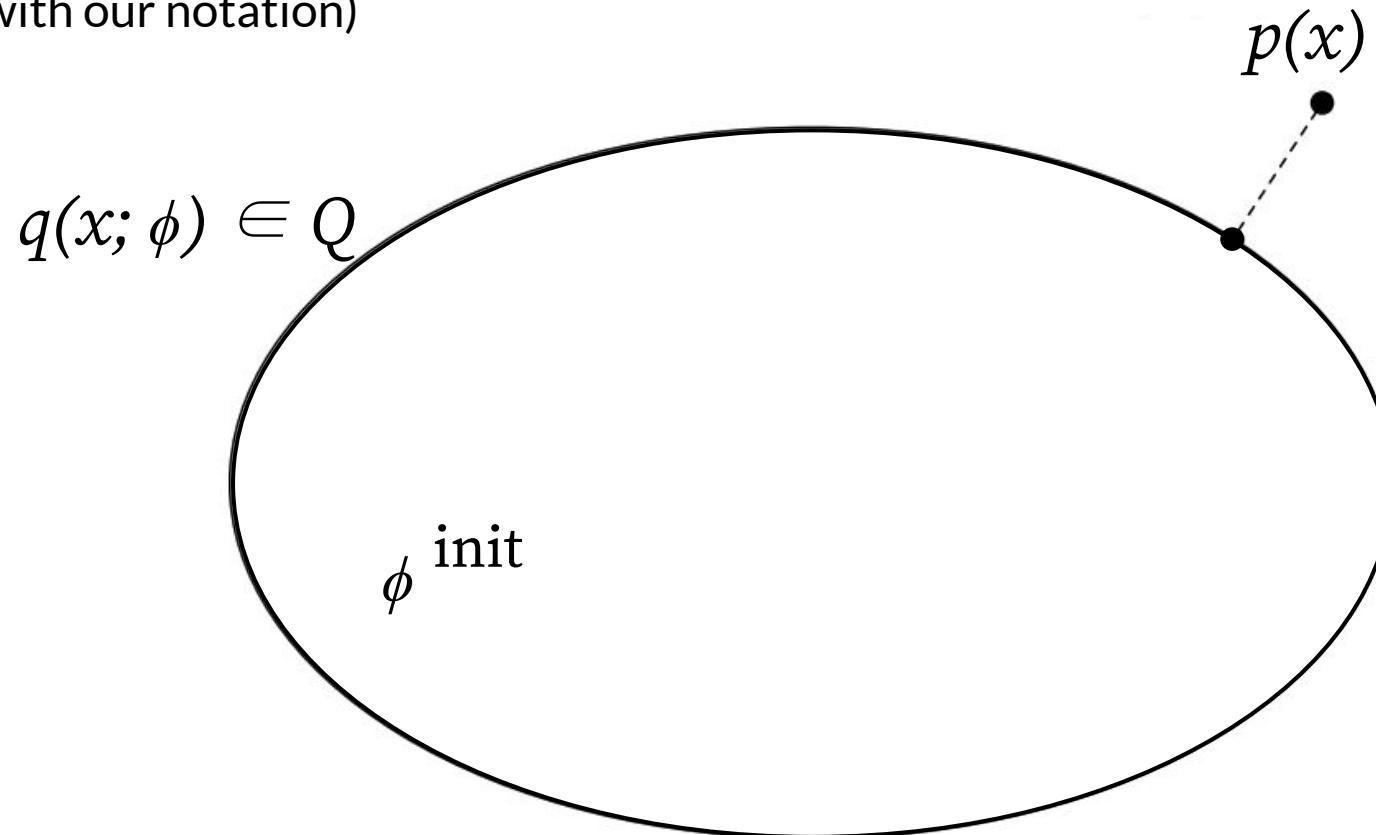
Variational Inference – Overview

An illustration from Dave Blei's Variational Inference Tutorial:
(replaced with our notation)



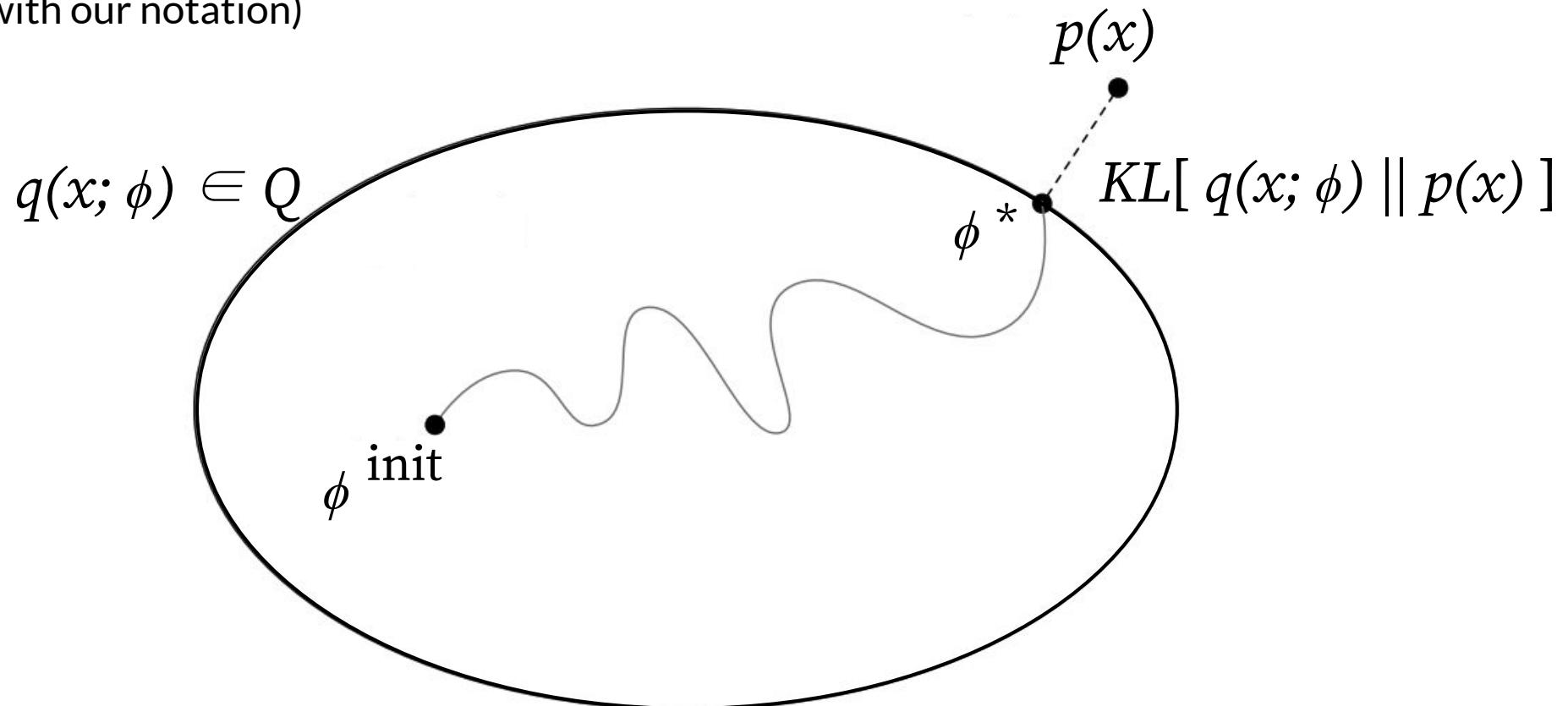
Variational Inference – Overview

An illustration from Dave Blei's Variational Inference Tutorial:
(replaced with our notation)



Variational Inference – Overview

An illustration from Dave Blei's Variational Inference Tutorial:
(replaced with our notation)



Kullback-Leibler Divergence

More formally:

Kullback-Leibler Divergence

More formally:

The KL divergence between distributions q and p with *discrete support* is:

$$\text{KL} [q \parallel p] = \sum_x q(x) \log \left(\frac{q(x)}{p(x)} \right) = \mathbb{E}_{q(x)} \left[\log \left(\frac{q(x)}{p(x)} \right) \right]$$

Kullback-Leibler Divergence

More formally:

The KL divergence between distributions q and p with *discrete support* is:

$$\text{KL} [q \parallel p] = \sum_x q(x) \log \left(\frac{q(x)}{p(x)} \right) = \mathbb{E}_{q(x)} \left[\log \left(\frac{q(x)}{p(x)} \right) \right]$$

And the KL divergence between distributions q and p with *continuous support* is:

$$\text{KL} [q \parallel p] = \int_{x \in \mathcal{X}} q(x) \log \left(\frac{q(x)}{p(x)} \right) dx = \mathbb{E}_{q(x)} \left[\log \left(\frac{q(x)}{p(x)} \right) \right]$$

Kullback-Leibler Divergence

Intuitively, when minimizing KL divergence in variational inference:
(in the words of Dave Blei)

Kullback-Leibler Divergence

$$\text{KL} [q \parallel p] = \sum_x q(x) \log \left(\frac{q(x)}{p(x)} \right)$$

Intuitively, when minimizing KL divergence in variational inference:
(in the words of Dave Blei)

Kullback-Leibler Divergence

$$\text{KL} [q \parallel p] = \sum_x q(x) \log \left(\frac{q(x)}{p(x)} \right)$$

Intuitively, when minimizing KL divergence in variational inference:

(in the words of Dave Blei)

- If $q(x)$ is high and $p(x)$ is also high (or higher) then we are happy. 😊

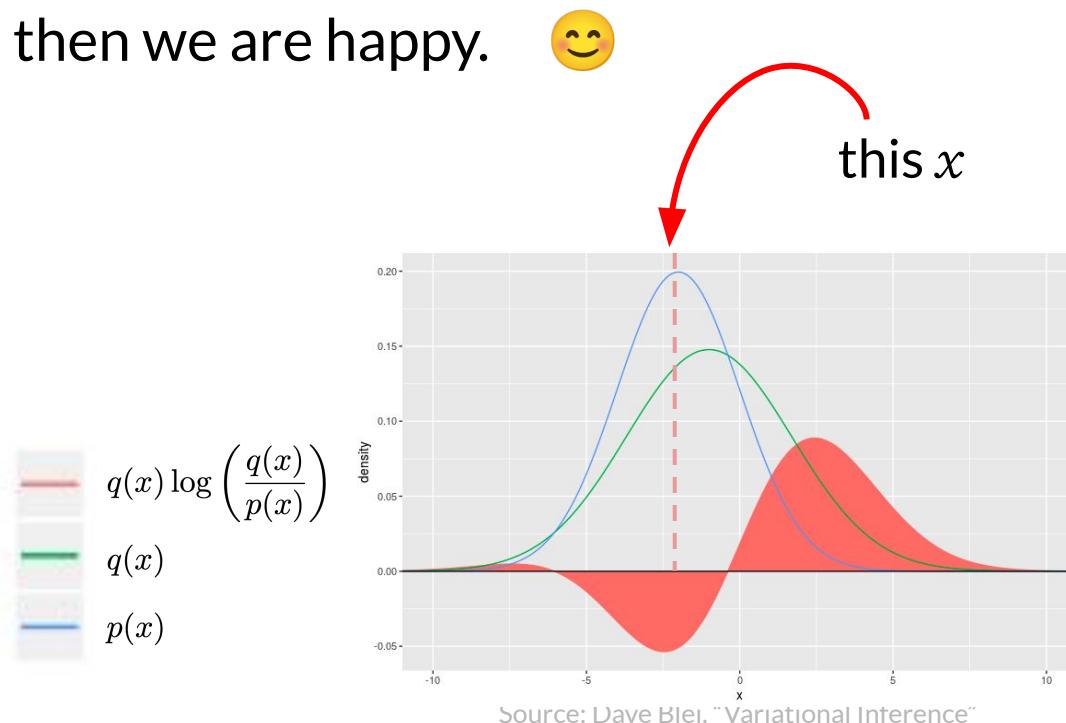
Kullback-Leibler Divergence

$$\text{KL} [q \parallel p] = \sum_x q(x) \log \left(\frac{q(x)}{p(x)} \right)$$

Intuitively, when minimizing KL divergence in variational inference:

(in the words of Dave Blei)

- If $q(x)$ is high and $p(x)$ is also high (or higher) then we are happy. 😊



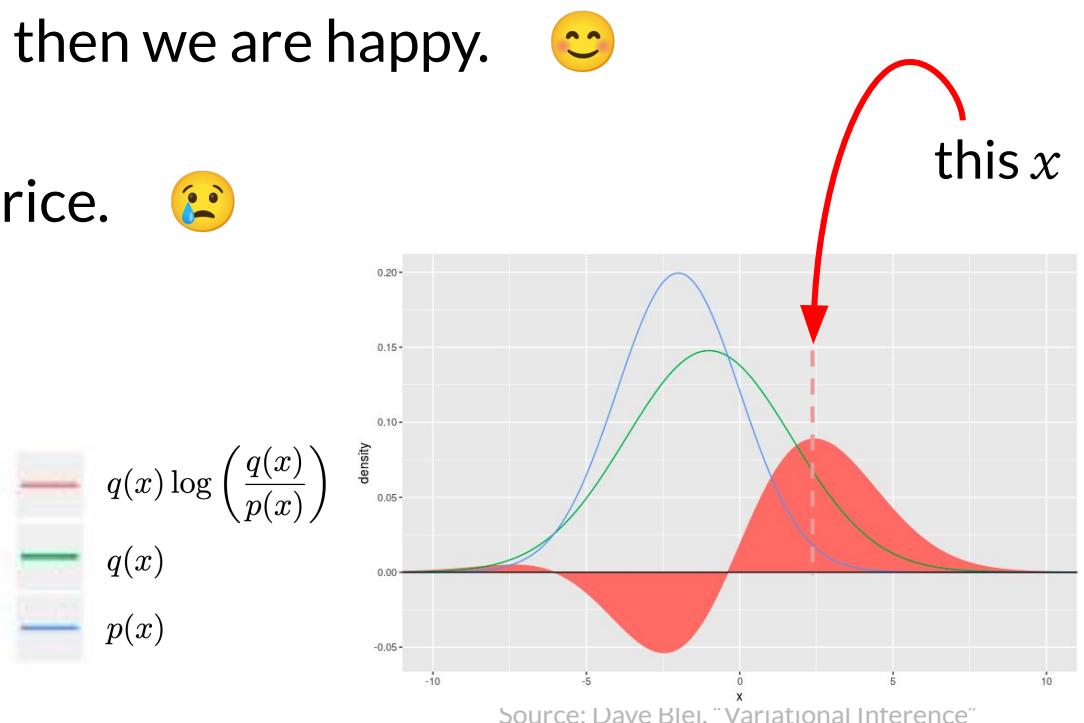
Kullback-Leibler Divergence

$$\text{KL} [q \parallel p] = \sum_x q(x) \log \left(\frac{q(x)}{p(x)} \right)$$

Intuitively, when minimizing KL divergence in variational inference:

(in the words of Dave Blei)

- If $q(x)$ is high and $p(x)$ is also high (or higher) then we are happy. 😊
- If $q(x)$ is high and $p(x)$ is low, then we pay a price. 😢

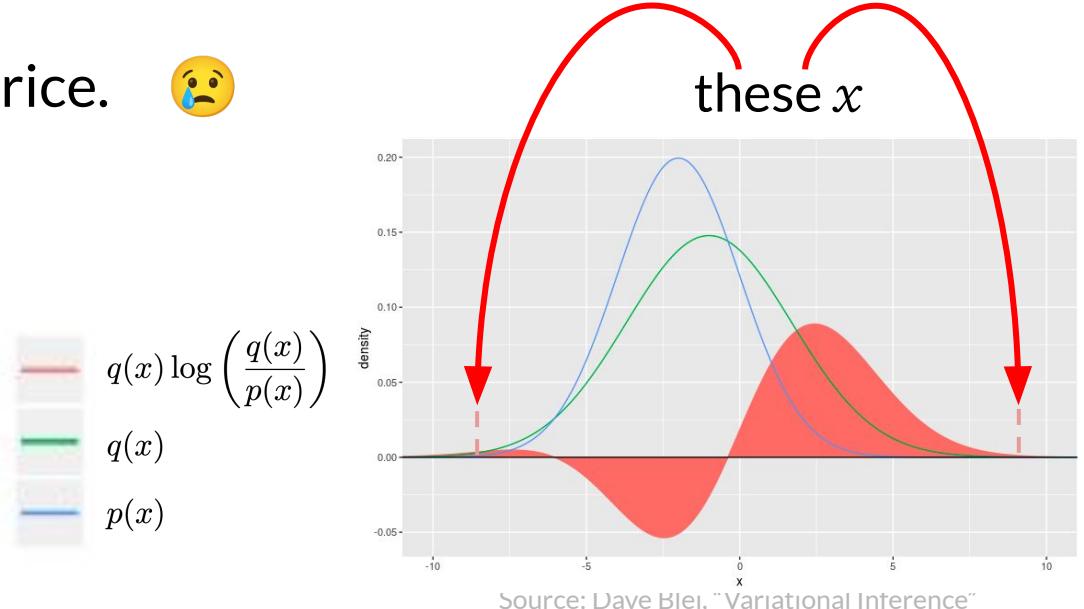
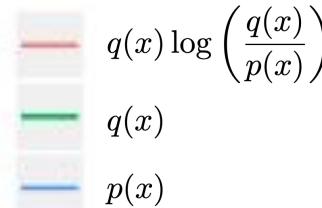


Kullback-Leibler Divergence

$$\text{KL} [q \parallel p] = \sum_x q(x) \log \left(\frac{q(x)}{p(x)} \right)$$

Intuitively, when minimizing KL divergence in variational inference:
(in the words of Dave Blei)

- If $q(x)$ is high and $p(x)$ is also high (or higher) then we are happy. 😊
- If $q(x)$ is high and $p(x)$ is low, then we pay a price. 😢
- If $q(x)$ is low... then we don't care. 😐



The Evidence Lower Bound

The Evidence Lower Bound

Suppose we are given a probability model of the form: $p_{\theta}(x) = \frac{1}{Z(\theta)} \prod_k \phi_k(x_k; \theta)$

The Evidence Lower Bound

A general way of writing
many PDFs of interest

Suppose we are given a probability model of the form: $p_{\theta}(x) = \frac{1}{Z(\theta)} \prod_k \phi_k(x_k; \theta)$

The Evidence Lower Bound

Suppose we are given a probability model of the form: $p_{\theta}(x) = \frac{1}{Z(\theta)} \prod_k \phi_k(x_k; \theta)$

And that we want to perform variational inference (*i.e.*, optimize $\text{KL}[q \parallel p]$).

The Evidence Lower Bound

Suppose we are given a probability model of the form: $p_{\theta}(x) = \frac{1}{Z(\theta)} \prod_k \phi_k(x_k; \theta)$

And that we want to perform variational inference (*i.e.*, optimize $\text{KL}[q \parallel p]$).

It is difficult to optimize $\text{KL}[q \parallel p]$ directly! For a few reasons:

The Evidence Lower Bound

Suppose we are given a probability model of the form: $p_\theta(x) = \frac{1}{Z(\theta)} \prod_k \phi_k(x_k; \theta)$

And that we want to perform variational inference (i.e., optimize $\text{KL}[q \parallel p]$).

It is difficult to optimize $\text{KL}[q \parallel p]$ directly! For a few reasons:

- (Often) intractable normalization constant.
- \Rightarrow difficult to take gradients due to this.
- In fact, it's difficult even to evaluate $\text{KL}[q \parallel p]$ (or $p(x)$ for that matter).

The Evidence Lower Bound

Therefore, we will not work with the KL divergence directly.

The Evidence Lower Bound

Therefore, we will not work with the KL divergence directly.

Instead, we will use the following objective, which has a similar form as the KL divergence, but only involves the *unnormalized probability* $\tilde{p}(x) = \prod_k \phi_k(x_k; \theta)$.

The Evidence Lower Bound

Therefore, we will not work with the KL divergence directly.

Instead, we will use the following objective, which has a similar form as the KL divergence, but only involves the *unnormalized probability* $\tilde{p}(x) = \prod_k \phi_k(x_k; \theta)$.

This objective is:

The Evidence Lower Bound

Therefore, we will not work with the KL divergence directly.

Instead, we will use the following objective, which has a similar form as the KL divergence, but only involves the *unnormalized probability* $\tilde{p}(x) = \prod_k \phi_k(x_k; \theta)$.

This objective is:

$$J(q) = \sum_x q(x) \log \frac{q(x)}{\tilde{p}(x)} = \mathbb{E}_{q(x)} \left[\log \frac{q(x)}{\tilde{p}(x)} \right]$$

For discrete distributions

The Evidence Lower Bound

Therefore, we will not work with the KL divergence directly.

Instead, we will use the following objective, which has a similar form as the KL divergence, but only involves the *unnormalized probability* $\tilde{p}(x) = \prod_k \phi_k(x_k; \theta)$.

This objective is:

$$J(q) = \sum_x q(x) \log \frac{q(x)}{\tilde{p}(x)} = \mathbb{E}_{q(x)} \left[\log \frac{q(x)}{\tilde{p}(x)} \right]$$

For discrete distributions

$$J(q) = \int_{\mathcal{X}} q(x) \log \frac{q(x)}{\tilde{p}(x)} dx = \mathbb{E}_{q(x)} \left[\log \frac{q(x)}{\tilde{p}(x)} \right]$$

For continuous distributions

The Evidence Lower Bound

This function is *tractable* (we can evaluate it), and it has the following property:

The Evidence Lower Bound

This function is *tractable* (we can evaluate it), and it has the following property:

$$\begin{aligned} J(q) &= \sum_x q(x) \log \frac{q(x)}{\tilde{p}(x)} \\ &= \sum_x q(x) \log \frac{q(x)}{p(x)} - \log Z(\theta) \\ &= \text{KL}[q(x) \parallel p(x)] - \log Z(\theta) \end{aligned}$$

Assuming a discrete distribution WLOG

The Evidence Lower Bound

This function is *tractable* (we can evaluate it), and it has the following property:

$$J(q) = \sum_x q(x) \log \frac{q(x)}{\tilde{p}(x)}$$

Definition of $J(q)$

Assuming a discrete distribution WLOG

$$\begin{aligned} &= \sum_x q(x) \log \frac{q(x)}{p(x)} - \log Z(\theta) \\ &= \text{KL}[q(x) \parallel p(x)] - \log Z(\theta) \end{aligned}$$

The Evidence Lower Bound

This function is *tractable* (we can evaluate it), and it has the following property:

$$J(q) = \sum_x q(x) \log \frac{q(x)}{\tilde{p}(x)}$$

Definition of $J(q)$

Assuming a discrete distribution WLOG

$$= \sum_x q(x) \log \frac{q(x)}{p(x)} - \log Z(\theta)$$

Definition of $p(x)$

$$= \text{KL} [q(x) \parallel p(x)] - \log Z(\theta)$$

The Evidence Lower Bound

This function is *tractable* (we can evaluate it), and it has the following property:

$$J(q) = \sum_x q(x) \log \frac{q(x)}{\tilde{p}(x)}$$

Definition of $J(q)$

Assuming a discrete distribution WLOG

$$= \sum_x q(x) \log \frac{q(x)}{p(x)} - \log Z(\theta)$$

Definition of $p(x)$

$$= \text{KL} [q(x) \parallel p(x)] - \log Z(\theta)$$

Definition of KL

The Evidence Lower Bound

This function is *tractable* (we can evaluate it), and it has the following property:

$$J(q) = \sum_x q(x) \log \frac{q(x)}{\tilde{p}(x)}$$

Definition of $J(q)$

Assuming a discrete distribution WLOG

$$= \sum_x q(x) \log \frac{q(x)}{p(x)} - \log Z(\theta)$$

Definition of $p(x)$

$$= \text{KL} [q(x) \parallel p(x)] - \log Z(\theta)$$

Definition of KL

$\Rightarrow J(q)$ is equal to the KL divergence minus the *log partition function* (aka *log normalizer*).
(aka “marginal log-likelihood”).

The Evidence Lower Bound

Also note that KL divergence is non-negative: $\text{KL} [q \parallel p] \geq 0$

Which implies the following relation:

The Evidence Lower Bound

Also note that KL divergence is non-negative: $\text{KL} [q \parallel p] \geq 0$

Which implies the following relation:

$$\log Z(\theta) = \text{KL} [q(x) \parallel p(x)] - J(q)$$

The Evidence Lower Bound

Also note that KL divergence is non-negative: $\text{KL} [q \parallel p] \geq 0$

Which implies the following relation:

Since KL is 0, at lowest

$$\log Z(\theta) = \text{KL} [q(x) \parallel p(x)] - J(q) \geq -J(q)$$

The Evidence Lower Bound

Also note that KL divergence is non-negative: $\text{KL} [q \parallel p] \geq 0$

Which implies the following relation:

Since KL is 0, at lowest

$$\log Z(\theta) = \text{KL} [q(x) \parallel p(x)] - J(q) \geq -J(q)$$

$$\Rightarrow \log Z(\theta) \geq -J(q)$$

The Evidence Lower Bound

Also note that KL divergence is non-negative: $\text{KL} [q \parallel p] \geq 0$

Which implies the following relation:

Since KL is 0, at lowest

$$\log Z(\theta) = \text{KL} [q(x) \parallel p(x)] - J(q) \geq -J(q)$$

$$\Rightarrow \log Z(\theta) \geq -J(q)$$

Therefore, $-J(q)$ is a *lower bound* on the log-partition function, $\log Z(\theta)$.



This is an important quantity! We call $-J(q)$ the **ELBO**.

The Evidence Lower Bound – Why do we call it the ELBO?

The Evidence Lower Bound – Why do we call it the ELBO?

Recall that we have been writing probabilistic models over $x = (x_1, \dots, x_n)$ as a normalized product of factors:

$$p_\theta(x_1, \dots, x_n) = \frac{1}{Z(\theta)} \tilde{p}_\theta(x_1, \dots, x_n) = \frac{1}{Z(\theta)} \prod_k \phi_k(x_k; \theta)$$

The Evidence Lower Bound – Why do we call it the ELBO?

Recall that we have been writing probabilistic models over $x = (x_1, \dots, x_n)$ as a normalized product of factors:

$$p_\theta(x_1, \dots, x_n) = \frac{1}{Z(\theta)} \tilde{p}_\theta(x_1, \dots, x_n) = \frac{1}{Z(\theta)} \prod_k \phi_k(x_k; \theta)$$

A distribution
(PDF) of interest

The Evidence Lower Bound – Why do we call it the ELBO?

Recall that we have been writing probabilistic models over $x = (x_1, \dots, x_n)$ as a normalized product of factors:

$$p_{\theta}(x_1, \dots, x_n) = \frac{1}{Z(\theta)} \tilde{p}_{\theta}(x_1, \dots, x_n) = \frac{1}{Z(\theta)} \prod_k \phi_k(x_k; \theta)$$

A distribution
(PDF) of interest

Unnormalized PDF
with normalization

The Evidence Lower Bound – Why do we call it the ELBO?

Recall that we have been writing probabilistic models over $x = (x_1, \dots, x_n)$ as a normalized product of factors:

$$p_{\theta}(x_1, \dots, x_n) = \frac{1}{Z(\theta)} \tilde{p}_{\theta}(x_1, \dots, x_n) = \frac{1}{Z(\theta)} \prod_k \phi_k(x_k; \theta)$$

A distribution
(PDF) of interest

Unnormalized PDF
with normalization

Normalized product
of factors

The Evidence Lower Bound – Why do we call it the ELBO?

As an example

The Evidence Lower Bound – Why do we call it the ELBO?

As an example, consider **posterior distributions** (of latents given observed), *i.e.*, $p(y \mid x_1, \dots, x_n)$. We can write

The Evidence Lower Bound – Why do we call it the ELBO?

As an example, consider **posterior distributions** (of latents given observed), i.e., $p(y \mid x_1, \dots, x_n)$. We can write

$$p_\theta(y \mid x_1, \dots, x_n) = \frac{p_\theta(y, x_1, \dots, x_n)}{p_\theta(x_1, \dots, x_n)} = \frac{1}{Z(\theta)} p_\theta(y, x_1, \dots, x_n)$$

The Evidence Lower Bound – Why do we call it the ELBO?

As an example, consider **posterior distributions** (of latents given observed), i.e., $p(y \mid x_1, \dots, x_n)$. We can write

$$p_\theta(y \mid x_1, \dots, x_n) = \frac{p_\theta(y, x_1, \dots, x_n)}{p_\theta(x_1, \dots, x_n)} = \frac{1}{Z(\theta)} p_\theta(y, x_1, \dots, x_n)$$

Definition of
conditional PDF

The Evidence Lower Bound – Why do we call it the ELBO?

As an example, consider **posterior distributions** (of latents given observed), i.e., $p(y \mid x_1, \dots, x_n)$. We can write

$$p_\theta(y \mid x_1, \dots, x_n) = \frac{p_\theta(y, x_1, \dots, x_n)}{p_\theta(x_1, \dots, x_n)} = \frac{1}{Z(\theta)} p_\theta(y, x_1, \dots, x_n)$$

Definition of
conditional PDF

Write as normalized
joint PDF

The Evidence Lower Bound – Why do we call it the ELBO?

As an example, consider **posterior distributions** (of latents given observed), i.e., $p(y \mid x_1, \dots, x_n)$. We can write

$$p_\theta(y \mid x_1, \dots, x_n) = \frac{p_\theta(y, x_1, \dots, x_n)}{p_\theta(x_1, \dots, x_n)} = \frac{1}{Z(\theta)} p_\theta(y, x_1, \dots, x_n)$$

Here, the partition function is $Z(\theta) = p_\theta(x_1, \dots, x_n)$, often called the **evidence**.
(or “marginal likelihood”)

The Evidence Lower Bound – Why do we call it the ELBO?

As an example, consider **posterior distributions** (of latents given observed), i.e., $p(y \mid x_1, \dots, x_n)$. We can write

$$p_\theta(y \mid x_1, \dots, x_n) = \frac{p_\theta(y, x_1, \dots, x_n)}{p_\theta(x_1, \dots, x_n)} = \frac{1}{Z(\theta)} p_\theta(y, x_1, \dots, x_n)$$

Here, the partition function is $Z(\theta) = p_\theta(x_1, \dots, x_n)$, often called the **evidence**.
(or “marginal likelihood”)

Thus, $-J(q) \leq \log Z(\theta)$ is often referred to as the **evidence lower bound**, or **ELBO**.

The Evidence Lower Bound

To summarize, the ELBO is:

The Evidence Lower Bound

To summarize, the ELBO is:

For discrete distributions

$$-J(q) = - \sum_x q(x) \log \frac{q(x)}{\tilde{p}(x)} = \sum_x q(x) \log \frac{\tilde{p}(x)}{q(x)} = \mathbb{E}_{q(x)} \left[\log \frac{\tilde{p}(x)}{q(x)} \right]$$

The Evidence Lower Bound

To summarize, the ELBO is:

For discrete distributions

$$-J(q) = - \sum_x q(x) \log \frac{q(x)}{\tilde{p}(x)} = \sum_x q(x) \log \frac{\tilde{p}(x)}{q(x)} = \mathbb{E}_{q(x)} \left[\log \frac{\tilde{p}(x)}{q(x)} \right]$$

For continuous distributions

$$-J(q) = - \int_{\mathcal{X}} q(x) \log \frac{q(x)}{\tilde{p}(x)} dx = \int_{\mathcal{X}} q(x) \log \frac{\tilde{p}(x)}{q(x)} dx = \mathbb{E}_{q(x)} \left[\log \frac{\tilde{p}(x)}{q(x)} \right]$$

The Evidence Lower Bound

Also, the ELBO is often written in the form:

$$\begin{aligned} -J(q) = \mathbb{E}_{q(x)} [\log \tilde{p}(x) - \log q(x)] &= \log Z(\theta) - \text{KL} [q(x) \parallel p(x)] \\ &\leq \log Z(\theta) \end{aligned}$$

The Evidence Lower Bound

Also, the ELBO is often written in the form:

$$\begin{aligned} -J(q) = \mathbb{E}_{q(x)} [\log \tilde{p}(x) - \log q(x)] &= \log Z(\theta) - \text{KL}[q(x) \parallel p(x)] \\ &\leq \log Z(\theta) \end{aligned}$$

The Evidence Lower Bound

Also, the ELBO is often written in the form:

$$\begin{aligned} -J(q) = \mathbb{E}_{q(x)} [\log \tilde{p}(x) - \log q(x)] &= \log Z(\theta) - \text{KL} [q(x) \parallel p(x)] \\ &\leq \log Z(\theta) \end{aligned}$$

The Evidence Lower Bound

Also, the ELBO is often written in the form:

$$\begin{aligned} -J(q) = \mathbb{E}_{q(x)} [\log \tilde{p}(x) - \log q(x)] &= \log Z(\theta) - \text{KL} [q(x) \parallel p(x)] \\ &\leq \log Z(\theta) \end{aligned}$$

The Evidence Lower Bound

Also, the ELBO is often written in the form:

$$\begin{aligned} -J(q) = \mathbb{E}_{q(x)} [\log \tilde{p}(x) - \log q(x)] &= \log Z(\theta) - \text{KL} [q(x) \parallel p(x)] \\ &\leq \log Z(\theta) \end{aligned}$$

Importantly, note that:

The Evidence Lower Bound

Also, the ELBO is often written in the form:

$$\begin{aligned} -J(q) = \mathbb{E}_{q(x)} [\log \tilde{p}(x) - \log q(x)] &= \log Z(\theta) - \text{KL}[q(x) \parallel p(x)] \\ &\leq \log Z(\theta) \end{aligned}$$

Importantly, note that:

- The difference between $\log Z(\theta)$ and $-J(q)$ is $\text{KL}[q \parallel p]$.

The Evidence Lower Bound

Also, the ELBO is often written in the form:

$$\begin{aligned} -J(q) = \mathbb{E}_{q(x)} [\log \tilde{p}(x) - \log q(x)] &= \log Z(\theta) - \text{KL} [q(x) \parallel p(x)] \\ &\leq \log Z(\theta) \end{aligned}$$

Importantly, note that:

- The difference between $\log Z(\theta)$ and $-J(q)$ is $\text{KL} [q \parallel p]$.
- Therefore, by maximizing the ELBO, we are *minimizing* $\text{KL} [q \parallel p]...$

The Evidence Lower Bound

Also, the ELBO is often written in the form:

$$\begin{aligned} -J(q) = \mathbb{E}_{q(x)} [\log \tilde{p}(x) - \log q(x)] &= \log Z(\theta) - \text{KL} [q(x) \parallel p(x)] \\ &\leq \log Z(\theta) \end{aligned}$$

Importantly, note that:

- The difference between $\log Z(\theta)$ and $-J(q)$ is $\text{KL} [q \parallel p]$.
- Therefore, by maximizing the ELBO, we are *minimizing* $\text{KL} [q \parallel p]...$
- ... by “*squeezing*” it between $-J(q)$ and $\log Z(\theta)$.

Variational Autoencoder (VAE)

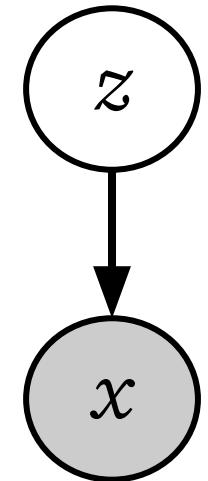
Setup

Setup

Example PGM

For simplicity, consider a latent variable model of the form:

$$p_{\theta}(x, z) = p_{\theta}(x \mid z)p_{\theta}(z)$$



Setup

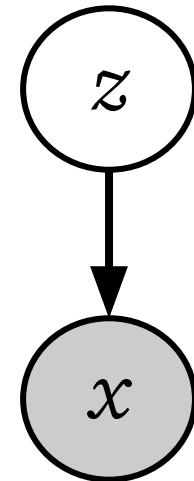
Example PGM

For simplicity, consider a latent variable model of the form:

$$p_{\theta}(x, z) = p_{\theta}(x \mid z)p_{\theta}(z)$$

where:

- $x \in \mathcal{X}$ are observed variables.
 - (\mathcal{X} can be either discrete or continuous.)
- $z \in \mathbb{R}^k$ are latent variables.



Setup

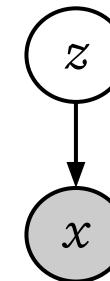
As an example, for intuition, you can imagine:

Setup

As an example, for intuition, you can imagine:

- x to be an image (e.g., of a human face).

Example PGM

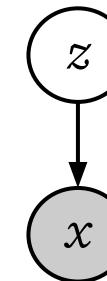


Setup

As an example, for intuition, you can imagine:

- x to be an image (e.g., of a human face).
- z to be latent factors, which explain features of the face, e.g.:
 - One coordinate could encode whether the face is happy or sad.
 - One coordinate could encode whether the face has short vs long hair.
 - One coordinate could encode whether the face has a beard vs clean shaven.
 - etc.

Example PGM



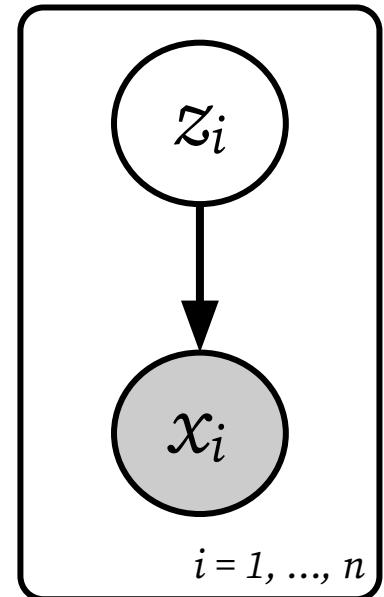
Setup

Suppose now that we have a dataset: $D = \{x_1, x_2, \dots, x_n\}$

Setup

Suppose now that we have a dataset: $D = \{x_1, x_2, \dots, x_n\}$

Example PGM

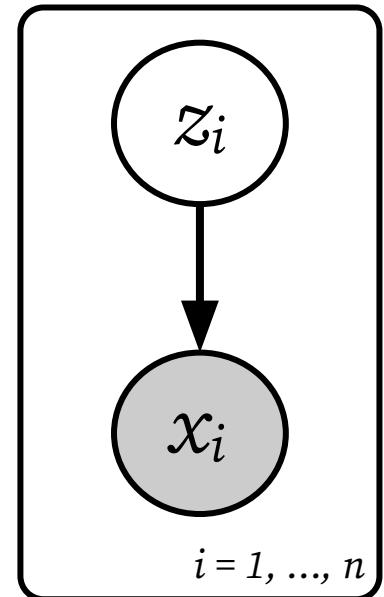


Setup

Suppose now that we have a dataset: $D = \{x_1, x_2, \dots, x_n\}$

We are interested in the following inference and learning tasks:

Example PGM



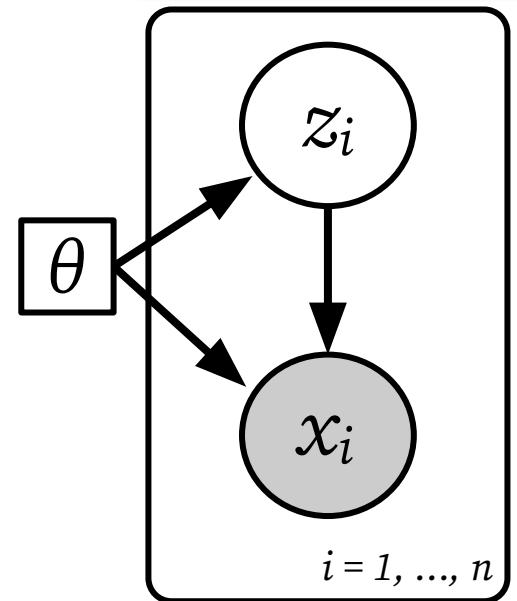
Setup

Suppose now that we have a dataset: $D = \{x_1, x_2, \dots, x_n\}$

We are interested in the following inference and learning tasks:

- (1) Learning the parameters θ of our model p_θ .

Example PGM



Setup

Suppose now that we have a dataset: $D = \{x_1, x_2, \dots, x_n\}$

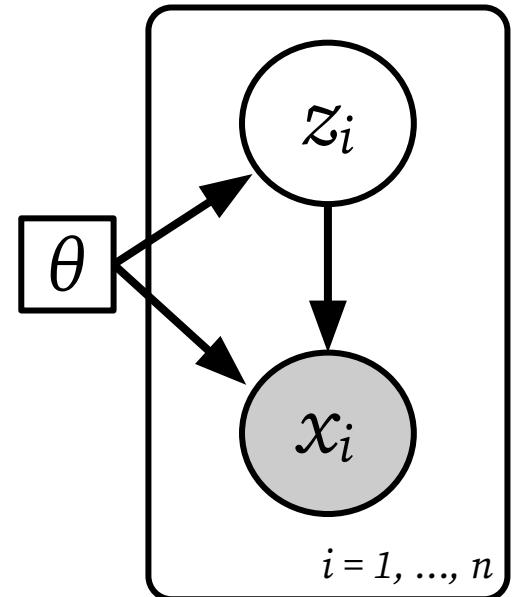
We are interested in the following inference and learning tasks:

(1) Learning the parameters θ of our model p_θ .

(2) Approximate posterior inference over z :

Given an image x , what are its latent factors, i.e., what is $p_\theta(z | x)$?

Example PGM



Background: VAE Paper

(From ICLR 2014)

Auto-Encoding Variational Bayes

Diederik P. Kingma
Machine Learning Group
Universiteit van Amsterdam
dpkingma@gmail.com

Max Welling
Machine Learning Group
Universiteit van Amsterdam
welling.max@gmail.com

Abstract

How can we perform efficient inference and learning in directed probabilistic models, in the presence of continuous latent variables with intractable posterior distributions, and large datasets? We introduce a stochastic variational inference and learning algorithm that scales to large datasets and, under some mild differentiability conditions, even works in the intractable case. Our contributions is two-fold. First, we show that a reparameterization of the variational lower bound with an independent noise variable yields a lower bound estimator that can be jointly optimized w.r.t. variational and generative parameters using standard gradient-based stochastic optimization methods. Second, we show that posterior inference can be made especially efficient by optimizing a probabilistic encoder (also called a recognition model) to approximate the intractable posterior, using the proposed estimator. Theoretical advantages are reflected in experimental results.

1 Introduction

How can we efficiently learn the parameters of directed probabilistic models whose continuous latent variables have intractable posterior distributions? The variational approach to Bayesian inference involves the introduction of an approximation to the intractable posterior, used to maximize the variational lower bound on the marginal likelihood. Unfortunately, the common mean-field approach requires analytical solutions of expectations w.r.t. the approximate posterior, which are also intractable in the general case. We show how a reparameterization of the variational lower bound yields a practical differentiable estimator of the lower bound. This SGVB (Stochastic Gradient Variational Bayes) estimator can be straightforwardly used as a stochastic objective function, and that can be jointly optimized w.r.t. both the variational and generative parameters, using standard stochastic gradient ascent techniques.

The SGVB algorithm can be applied to learning almost any generative model with continuous latent variables. When we use a neural network for the posterior approximation, we arrive at a *variational auto-encoder*. The SGVB objective for this case contains a regularization term dictated by the variational bound, and a stochastic data reconstruction term. From the learned generative model it is straightforward to generate samples simply by ancestral sampling. The probabilistic encoder can be used for fast approximate inference of latent variables, i.e. for recognition, representation or visualization purposes. Furthermore, the lower bound estimator can be used for inference tasks such as denoising and inpainting.

2 Method

The strategy in this section can be used to derive a lower bound estimator (a stochastic objective function) for a variety of directed graphical models with continuous latent variables. We will restrict ourselves here to the common case where we have an i.i.d. dataset with latent variables per datapoint, and where we like to perform maximum likelihood (ML) or maximum a posteriori (MAP) inference

Background: VAE Paper

(From ICLR 2014)

Introduces two distinct things:

(1) Auto-encoding variational Bayes (AEVB) algorithm.
(a stochastic-gradient VI procedure)

(2) A special case (involving neural networks): the *variational autoencoder* (VAE).

Auto-Encoding Variational Bayes

Diederik P. Kingma
Machine Learning Group
Universiteit van Amsterdam
dpkingma@gmail.com

Max Welling
Machine Learning Group
Universiteit van Amsterdam
welling.max@gmail.com

Abstract

How can we perform efficient inference and learning in directed probabilistic models, in the presence of continuous latent variables with intractable posterior distributions, and large datasets? We introduce a stochastic variational inference and learning algorithm that scales to large datasets and, under some mild differentiability conditions, even works in the intractable case. Our contributions is two-fold. First, we show that a reparameterization of the variational lower bound with an independent noise variable yields a lower bound estimator that can be jointly optimized w.r.t. variational and generative parameters using standard gradient-based stochastic optimization methods. Second, we show that posterior inference can be made especially efficient by optimizing a probabilistic encoder (also called a recognition model) to approximate the intractable posterior, using the proposed estimator. Theoretical advantages are reflected in experimental results.

1 Introduction

How can we efficiently learn the parameters of directed probabilistic models whose continuous latent variables have intractable posterior distributions? The variational approach to Bayesian inference involves the introduction of an approximation to the intractable posterior, used to maximize the variational lower bound on the marginal likelihood. Unfortunately, the common mean-field approach requires analytical solutions of expectations w.r.t. the approximate posterior, which are also intractable in the general case. We show how a reparameterization of the variational lower bound yields a practical differentiable estimator of the lower bound. This SGVB (Stochastic Gradient Variational Bayes) estimator can be straightforwardly used as a stochastic objective function, and that can be jointly optimized w.r.t. both the variational and generative parameters, using standard stochastic gradient ascent techniques.

The SGVB algorithm can be applied to learning almost any generative model with continuous latent variables. When we use a neural network for the posterior approximation, we arrive at a *variational auto-encoder*. The SGVB objective for this case contains a regularization term dictated by the variational bound, and a stochastic data reconstruction term. From the learned generative model it is straightforward to generate samples simply by ancestral sampling. The probabilistic encoder can be used for fast approximate inference of latent variables, i.e. for recognition, representation or visualization purposes. Furthermore, the lower bound estimator can be used for inference tasks such as denoising and inpainting.

2 Method

The strategy in this section can be used to derive a lower bound estimator (a stochastic objective function) for a variety of directed graphical models with continuous latent variables. We will restrict ourselves here to the common case where we have an i.i.d. dataset with latent variables per datapoint, and where we like to perform maximum likelihood (ML) or maximum a posteriori (MAP) inference

Auto-encoding Variational Bayes (AEVB)

Auto-encoding Variational Bayes (AEVB)

The Auto-encoding Variational Bayes (AEVB) algorithm aims to solve the tasks from before, in particular:

Auto-encoding Variational Bayes (AEVB)

The Auto-encoding Variational Bayes (AEVB) algorithm aims to solve the tasks from before, in particular:

- (1) Learning the parameters θ of our model p_θ .
- (2) Approximate posterior inference over z : $p_\theta(z \mid x)$

Auto-encoding Variational Bayes (AEVB)

The Auto-encoding Variational Bayes (AEVB) algorithm aims to solve the tasks from before, in particular:

- (1) Learning the parameters θ of our model p_θ .
- (2) Approximate posterior inference over z : $p_\theta(z | x)$

The task of generative modeling

Auto-encoding Variational Bayes (AEVB)

The Auto-encoding Variational Bayes (AEVB) algorithm aims to solve the tasks from before, in particular:

- (1) Learning the parameters θ of our model p_θ .
- (2) Approximate posterior inference over z : $p_\theta(z | x)$

The task of generative modeling

We'll start with methods for approximate posterior inference, and then show how this can lead to learning the parameters (and thus to generative modeling!)

Auto-encoding Variational Bayes (AEVB) – Family of Approximations

Auto-encoding Variational Bayes (AEVB) – Family of Approximations

Suppose we have a family of approximations for our model posterior $p_\theta(z \mid x)$, written $q_\phi(z \mid x)$, which is parameterized by some parameter ϕ .

Auto-encoding Variational Bayes (AEVB) – Family of Approximations

Suppose we have a family of approximations for our model posterior $p_\theta(z | x)$, written $q_\phi(z | x)$, which is parameterized by some parameter ϕ .

A note on the family of approximations $q_\phi(z | x)$:

- We are trying to approximate a posterior distribution for any given observation x .
- So, if we want, we can define a family of approximations as a function of x .
- This effectively yields a different variational approximation $q_\phi(z)$ for each x .

Auto-encoding Variational Bayes (AEVB) – The ELBO

Suppose we have a family of approximations for our model posterior $p_\theta(z \mid x)$, written $q_\phi(z \mid x)$, which is parameterized by some parameter ϕ .

Auto-encoding Variational Bayes (AEVB) – The ELBO

Suppose we have a family of approximations for our model posterior $p_\theta(z \mid x)$, written $q_\phi(z \mid x)$, which is parameterized by some parameter ϕ .

Consider an ELBO for this model, written as:

$$\mathcal{L}(p_\theta, q_\phi) = \mathbb{E}_{q_\phi(z|x)} [\log p_\theta(x, z) - \log q_\phi(z \mid x)]$$

Auto-encoding Variational Bayes (AEVB) – The ELBO

Suppose we have a family of approximations for our model posterior $p_\theta(z \mid x)$, written $q_\phi(z \mid x)$, which is parameterized by some parameter ϕ .

Consider an ELBO for this model, written as:

This is $-J(q)$
from before!

$$\mathcal{L}(p_\theta, q_\phi) = \mathbb{E}_{q_\phi(z|x)} [\log p_\theta(x, z) - \log q_\phi(z \mid x)]$$

Log of unnormalized
PDF

Log of approximate
PDF

Auto-encoding Variational Bayes (AEVB) – The ELBO

Suppose we have a family of approximations for our model posterior $p_\theta(z \mid x)$, written $q_\phi(z \mid x)$, which is parameterized by some parameter ϕ .

Consider an ELBO for this model, written as:

This is $-J(q)$
from before!

$$\mathcal{L}(p_\theta, q_\phi) = \mathbb{E}_{q_\phi(z|x)} [\log p_\theta(x, z) - \log q_\phi(z \mid x)]$$

Note that the ELBO satisfies:

$$\log p_\theta(x) = \text{KL} [q_\phi(z \mid x) \parallel p_\theta(z \mid x)] + \mathcal{L}(p_\theta, q_\phi).$$

Evidence equals KL divergence plus ELBO.

Auto-encoding Variational Bayes (AEVB) – BBVI

How can we optimize the ELBO to find the optimal approximation $q_\phi(z | x)$?

Auto-encoding Variational Bayes (AEVB) – BBVI

How can we optimize the ELBO to find the optimal approximation $q_\phi(z | x)$?

⇒ We'll use a *black-box variational inference* approach.

This means a VI method that aims to be automatic / applied in a general purpose way.

Auto-encoding Variational Bayes (AEVB) – BBVI

How can we optimize the ELBO to find the optimal approximation $q_\phi(z | x)$?

⇒ We'll use a *black-box variational inference* approach.

This means a VI method that aims to be automatic / applied in a general purpose way.

In particular:

- Maximize the ELBO via gradient descent over ϕ .
- Simply requires q_ϕ to be differentiable with respect to ϕ .
- (Along with a couple of other things...)

Auto-encoding Variational Bayes (AEVB) – Learning and Inference

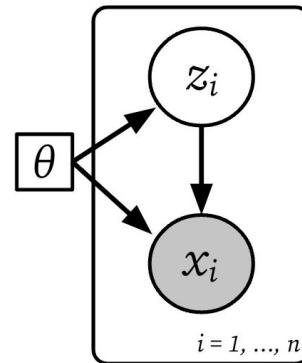
Furthermore

- Instead of just inference, we'll simultaneously do learning of parameters θ .

Auto-encoding Variational Bayes (AEVB) – Learning and Inference

Furthermore

- Instead of just inference, we'll simultaneously do learning of parameters θ .



Intuitively:

Simultaneously update our probabilistic model, *while* we are performing approximate posterior inference in the model.

Auto-encoding Variational Bayes (AEVB) – Learning and Inference

Furthermore

- Instead of just inference, we'll simultaneously do learning of parameters θ .
- \Rightarrow Via gradient descent over both ϕ and θ .

Auto-encoding Variational Bayes (AEVB) – Learning and Inference

Furthermore

- Instead of just inference, we'll simultaneously do learning of parameters θ .
- \Rightarrow Via gradient descent over both ϕ and θ .
- Optimization of ϕ moves the ELBO closer to $\log p(x)$ (i.e., minimizes KL divergence).

Auto-encoding Variational Bayes (AEVB) – Learning and Inference

Furthermore

- Instead of just inference, we'll simultaneously do learning of parameters θ .
- \Rightarrow Via gradient descent over both ϕ and θ .
- Optimization of ϕ moves the ELBO closer to $\log p(x)$ (i.e., minimizes KL divergence).
- Optimization of θ keeps pushing the upper bound (i.e., the evidence $\log p(x)$) up.
 - \Rightarrow Higher likelihood, better learning, and improved generative model.

Auto-encoding Variational Bayes (AEVB) – Learning and Inference

Furthermore

- Instead of just inference, we'll simultaneously do learning of parameters θ .
- \Rightarrow Via gradient descent over both ϕ and θ .
- Optimization of ϕ moves the ELBO closer to $\log p(x)$ (i.e., minimizes KL divergence).
- Optimization of θ keeps pushing the upper bound (i.e., the evidence $\log p(x)$) up.
 - \Rightarrow Higher likelihood, better learning, and improved generative model.

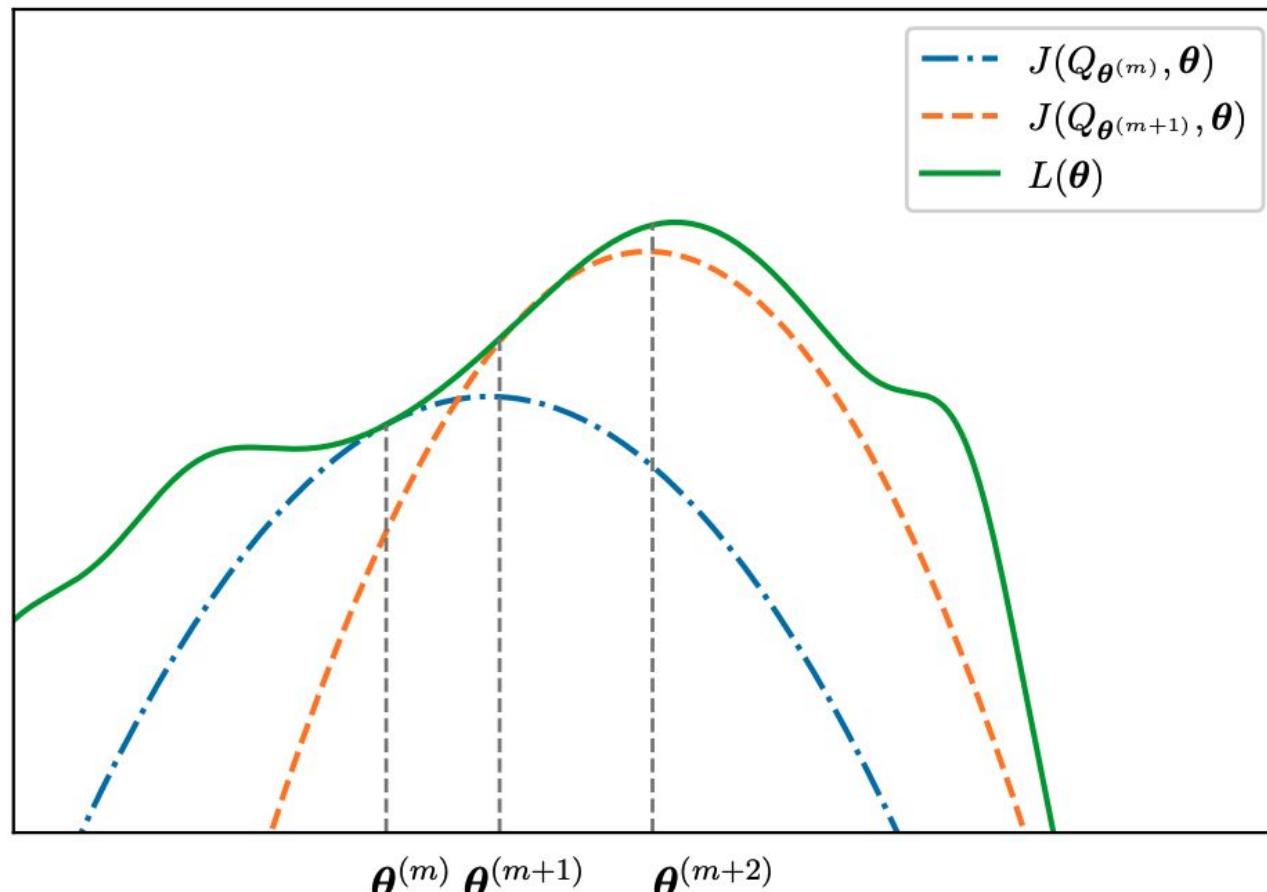
This is somewhat similar to the expectation-maximization (EM) algorithm...

Auto-encoding Variational Bayes (AEVB) – Learning and Inference

Recall the EM algorithm:

Auto-encoding Variational Bayes (AEVB) – Learning and Inference

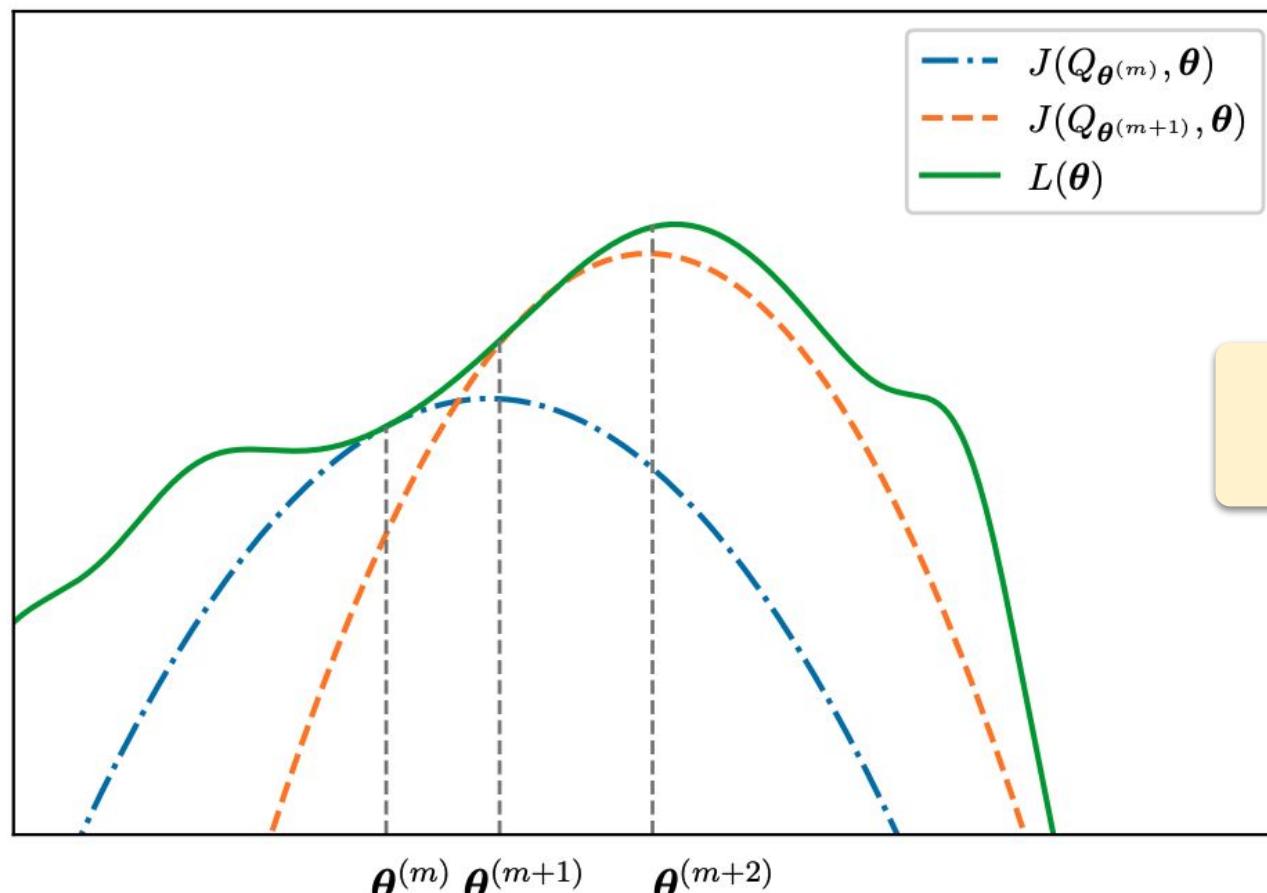
Recall the EM algorithm:



Iteratively forming then optimizing lower bounds to the (*marginal*) likelihood in latent-variable models.

Auto-encoding Variational Bayes (AEVB) – Learning and Inference

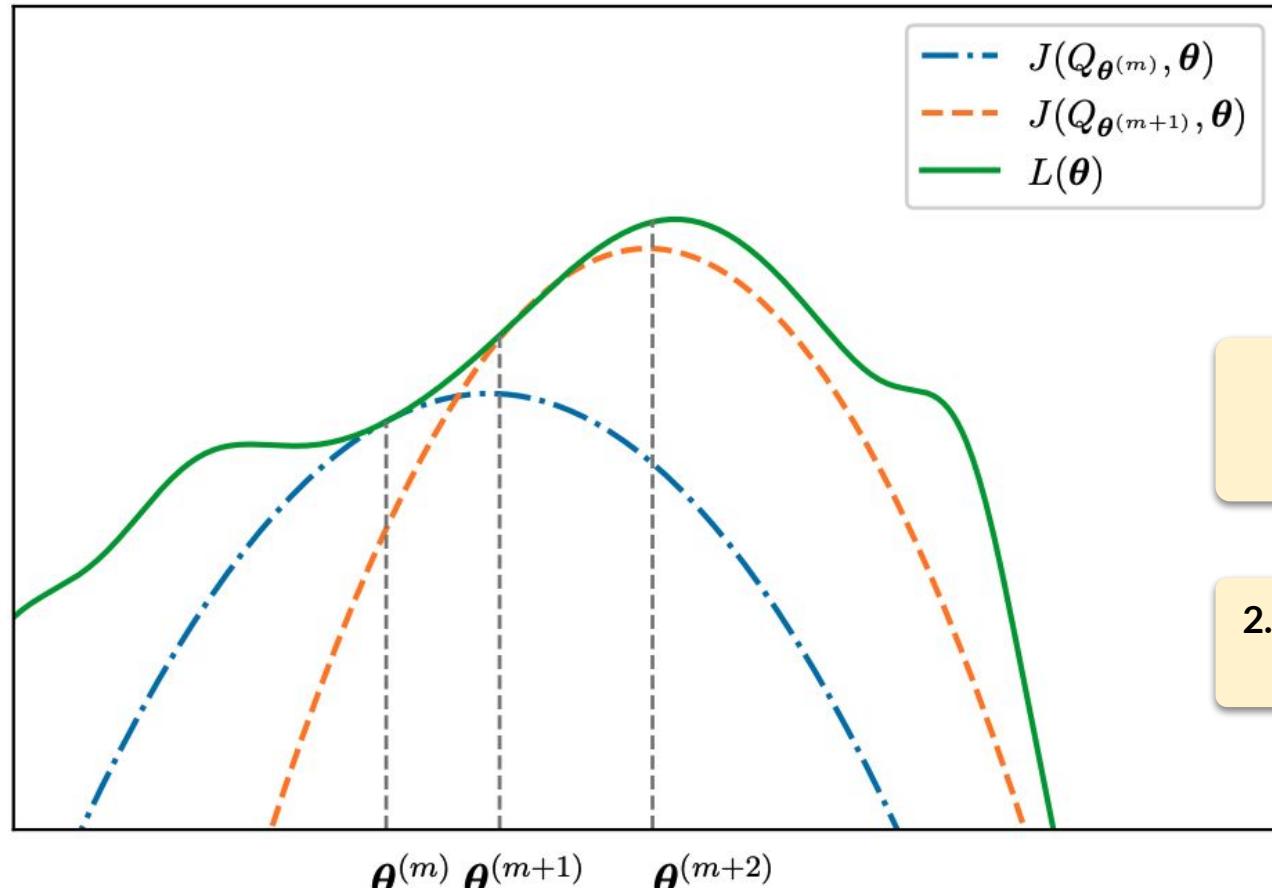
Recall the EM algorithm:



Iteratively forming then optimizing lower bounds to the (*marginal*) likelihood in latent-variable models.

Auto-encoding Variational Bayes (AEVB) – Learning and Inference

Recall the EM algorithm:



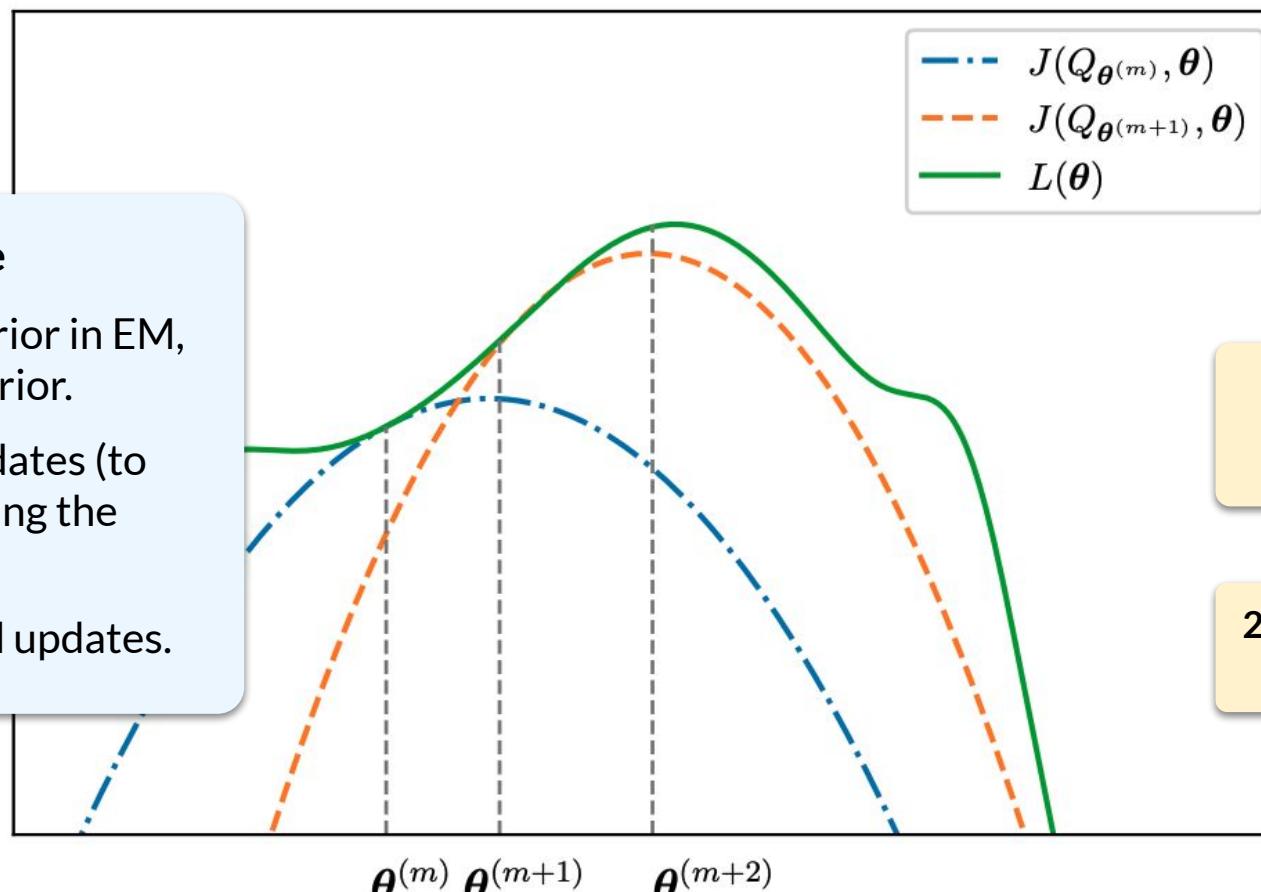
Iteratively forming then optimizing lower bounds to the (*marginal*) likelihood in latent-variable models.

1. E-Step: do inference to compute next objective (next lower bound)

2. M-Step: maximize this new objective (lower bound).

Auto-encoding Variational Bayes (AEVB) – Learning and Inference

Recall the EM algorithm:



Iteratively forming then optimizing lower bounds to the (*marginal*) likelihood in latent-variable models.

1. E-Step: do inference to compute next objective (next lower bound)

2. M-Step: maximize this new objective (lower bound).

Computing the Gradient

Computing the Gradient

So to optimize the ELBO, we need to compute the gradient:

$$\nabla_{\theta, \phi} \mathbb{E}_{q_\phi(z)} [\log p_\theta(x, z) - \log q_\phi(z)]$$

Computing the Gradient

So to optimize the ELBO, we need to compute the gradient:

Dropping the conditioning of q on x for now to keep it general.

$$\nabla_{\theta, \phi} \mathbb{E}_{q_\phi(z)} [\log p_\theta(x, z) - \log q_\phi(z)]$$



Computing the Gradient

So to optimize the ELBO, we need to compute the gradient:

Dropping the conditioning of q on x for now to keep it general.

$$\nabla_{\theta, \phi} \mathbb{E}_{q_\phi(z)} [\log p_\theta(x, z) - \log q_\phi(z)]$$

Note that an expression for this expectation in closed form is often *not possible*.

Computing the Gradient

So to optimize the ELBO, we need to compute the gradient:

Dropping the conditioning of q on x for now to keep it general.

$$\nabla_{\theta, \phi} \mathbb{E}_{q_\phi(z)} [\log p_\theta(x, z) - \log q_\phi(z)]$$

Note that an expression for this expectation in closed form is often *not possible*.



Instead, we could take gradient of a Monte Carlo estimate by sampling z from q_ϕ .

Computing the Gradient

So to optimize the ELBO, we need to compute the gradient:

Dropping the conditioning of q on x for now to keep it general.

$$\nabla_{\theta, \phi} \mathbb{E}_{q_\phi(z)} [\log p_\theta(x, z) - \log q_\phi(z)]$$

If we do this (gradient of MC estimate), note it is easy to write out the **gradient for p_θ** :

$$\begin{aligned} \nabla_{\theta} \mathbb{E}_{q_\phi(z)} [\log p_\theta(x, z) - \log q_\phi(z)] &= \mathbb{E}_{q_\phi(z)} [\nabla_{\theta} \log p_\theta(x, z)] \\ &\approx \frac{1}{M} \sum_{m=1}^M \nabla_{\theta} \log p_\theta(x, z^{(m)}) \\ &\quad \text{where } z^{(m)} \sim q_\phi(z) \end{aligned}$$

Computing the Gradient

So to optimize the ELBO, we need to compute the gradient:

Dropping the conditioning of q on x for now to keep it general.

$$\nabla_{\theta, \phi} \mathbb{E}_{q_\phi(z)} [\log p_\theta(x, z) - \log q_\phi(z)]$$

If we do this (gradient of MC estimate), note it is easy to write out the **gradient for p_θ** :

$$\begin{aligned} \nabla_{\theta} \mathbb{E}_{q_\phi(z)} [\log p_\theta(x, z) - \log q_\phi(z)] &= \mathbb{E}_{q_\phi(z)} [\nabla_{\theta} \log p_\theta(x, z)] \\ &\approx \frac{1}{M} \sum_{m=1}^M \nabla_{\theta} \log p_\theta(x, z^{(m)}) \end{aligned}$$

where $z^{(m)} \sim q_\phi(z)$

Computing the Gradient

So to optimize the ELBO, we need to compute the gradient:

Dropping the conditioning of q on x for now to keep it general.

$$\nabla_{\theta, \phi} \mathbb{E}_{q_\phi(z)} [\log p_\theta(x, z) - \log q_\phi(z)]$$

If we do this (gradient of MC estimate), note it is easy to write out the **gradient for p_θ** :

$$\nabla_{\theta} \mathbb{E}_{q_\phi(z)} [\log p_\theta(x, z) - \log q_\phi(z)] = \mathbb{E}_{q_\phi(z)} [\nabla_{\theta} \log p_\theta(x, z)]$$

$$\approx \frac{1}{M} \sum_{m=1}^M \nabla_{\theta} \log p_\theta(x, z^{(m)})$$

Just swap the gradient and expectation here.

where $z^{(m)} \sim q_\phi(z)$

Computing the Gradient

So to optimize the ELBO, we need to compute the gradient:

Dropping the conditioning of q on x for now to keep it general.

$$\nabla_{\theta, \phi} \mathbb{E}_{q_\phi(z)} [\log p_\theta(x, z) - \log q_\phi(z)]$$

If we do this (gradient of MC estimate), note it is easy to write out the **gradient for p_θ** :

$$\begin{aligned} \nabla_{\theta} \mathbb{E}_{q_\phi(z)} [\log p_\theta(x, z) - \log q_\phi(z)] &= \mathbb{E}_{q_\phi(z)} [\nabla_{\theta} \log p_\theta(x, z)] \\ &\approx \frac{1}{M} \sum_{m=1}^M \nabla_{\theta} \log p_\theta(x, z^{(m)}) \\ &\quad \text{where } z^{(m)} \sim q_\phi(z) \end{aligned}$$

Computing the Gradient

So to optimize the ELBO, we need to compute the gradient:

Dropping the conditioning of q on x for now to keep it general.

$$\nabla_{\theta, \phi} \mathbb{E}_{q_\phi(z)} [\log p_\theta(x, z) - \log q_\phi(z)]$$

But it's harder to write out the **gradient for q_ϕ** :

$$\nabla_\phi \mathbb{E}_{q_\phi(z)} [\log p_\theta(x, z) - \log q_\phi(z)] = ??$$

Computing the Gradient

So to optimize the ELBO, we need to compute the gradient:

Dropping the conditioning of q on x for now to keep it general.

$$\nabla_{\theta, \phi} \mathbb{E}_{q_\phi(z)} [\log p_\theta(x, z) - \log q_\phi(z)]$$

But it's harder to write out the **gradient for q_ϕ** :

$$\nabla_\phi \mathbb{E}_{q_\phi(z)} [\log p_\theta(x, z) - \log q_\phi(z)] = ??$$

Note that we cannot just swap the gradient and expectation, because the expectation is being taken with respect to $q_\phi(z)$.

The Score Function Gradient Estimator

One way to estimate this gradient is via the **score function estimator**:

The Score Function Gradient Estimator

One way to estimate this gradient is via the **score function estimator**:

$$\nabla_{\phi} \mathbb{E}_{q_{\phi}(z)} [\log p_{\theta}(x, z) - \log q_{\phi}(z)] =$$

$$\mathbb{E}_{q_{\phi}(z)} [(\log p_{\theta}(x, z) - \log q_{\phi}(z)) \nabla_{\phi} \log q_{\phi}(z)]$$

The Score Function Gradient Estimator

One way to estimate this gradient is via the **score function estimator**:

$$\nabla_{\phi} \mathbb{E}_{q_{\phi}(z)} [\log p_{\theta}(x, z) - \log q_{\phi}(z)] =$$

$$\mathbb{E}_{q_{\phi}(z)} [(\log p_{\theta}(x, z) - \log q_{\phi}(z)) \nabla_{\phi} \log q_{\phi}(z)]$$

Expectation with
respect to q of...

The Score Function Gradient Estimator

One way to estimate this gradient is via the **score function estimator**:

$$\nabla_{\phi} \mathbb{E}_{q_{\phi}(z)} [\log p_{\theta}(x, z) - \log q_{\phi}(z)] =$$

$$\mathbb{E}_{q_{\phi}(z)} [(\log p_{\theta}(x, z) - \log q_{\phi}(z)) \nabla_{\phi} \log q_{\phi}(z)]$$

Expectation with
respect to q of...

Log joint minus $\log q$
(same as in the ELBO)

The Score Function Gradient Estimator

One way to estimate this gradient is via the **score function estimator**:

$$\nabla_{\phi} \mathbb{E}_{q_{\phi}(z)} [\log p_{\theta}(x, z) - \log q_{\phi}(z)] =$$

$$\mathbb{E}_{q_{\phi}(z)} [(\log p_{\theta}(x, z) - \log q_{\phi}(z)) \nabla_{\phi} \log q_{\phi}(z)]$$

Expectation with
respect to q of...

Log joint minus log q
(same as in the ELBO)

Multiplied by the
gradient of the log of q .

The Score Function Gradient Estimator

One way to estimate this gradient is via the **score function estimator**:

$$\nabla_{\phi} \mathbb{E}_{q_{\phi}(z)} [\log p_{\theta}(x, z) - \log q_{\phi}(z)] =$$

$$\mathbb{E}_{q_{\phi}(z)} [(\log p_{\theta}(x, z) - \log q_{\phi}(z)) \nabla_{\phi} \log q_{\phi}(z)]$$

Expectation with
respect to q of...

Log joint minus log q
(same as in the ELBO)

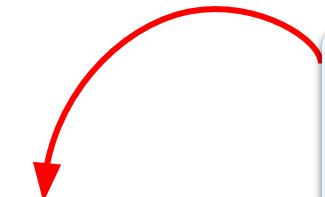
Multiplied by the
gradient of the log of q .

Follows from some basic algebra/calculus, and takes about half a page to derive :D.

The Score Function Gradient Estimator

One way to estimate this gradient is via the **score function estimator**:

$$\nabla_{\phi} \mathbb{E}_{q_{\phi}(z)} [\log p_{\theta}(x, z) - \log q_{\phi}(z)] =$$



Notably: since the gradient is now inside the expectation, we can evaluate this expectation using Monte Carlo methods!

$$\mathbb{E}_{q_{\phi}(z)} [(\log p_{\theta}(x, z) - \log q_{\phi}(z)) \nabla_{\phi} \log q_{\phi}(z)]$$

Expectation with
respect to q of...

Log joint minus log q
(same as in the ELBO)

Multiplied by the
gradient of the log of q .

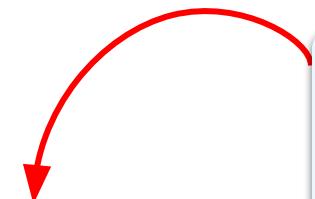
Follows from some basic algebra/calculus, and takes about half a page to derive :D.

⇒ Now can estimate via Monte Carlo methods.

The Score Function Gradient Estimator

One way to estimate this gradient is via the **score function estimator**:

$$\nabla_{\phi} \mathbb{E}_{q_{\phi}(z)} [\log p_{\theta}(x, z) - \log q_{\phi}(z)]$$



Notably: since the gradient is now inside the expectation, we can evaluate this expectation using Monte Carlo methods!

$$= \mathbb{E}_{q_{\phi}(z)} [(\log p_{\theta}(x, z) - \log q_{\phi}(z)) \nabla_{\phi} \log q_{\phi}(z)]$$

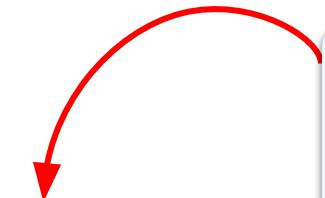
$$= \frac{1}{M} \sum_{m=1}^M \left(\log p_{\theta}(x, z^{(m)}) - \log q_{\phi}(z^{(m)}) \right) \nabla_{\phi} \log q_{\phi}(z^{(m)}),$$

where $z^{(m)} \sim q_{\phi}(z)$

The Score Function Gradient Estimator

One way to estimate this gradient is via the **score function estimator**:

$$\nabla_{\phi} \mathbb{E}_{q_{\phi}(z)} [\log p_{\theta}(x, z) - \log q_{\phi}(z)]$$



Notably: since the gradient is now inside the expectation, we can evaluate this expectation using Monte Carlo methods!

$$= \mathbb{E}_{q_{\phi}(z)} [(\log p_{\theta}(x, z) - \log q_{\phi}(z)) \nabla_{\phi} \log q_{\phi}(z)]$$

$$= \frac{1}{M} \sum_{m=1}^M \left(\log p_{\theta}(x, z^{(m)}) - \log q_{\phi}(z^{(m)}) \right) \nabla_{\phi} \log q_{\phi}(z^{(m)}),$$

where $z^{(m)} \sim q_{\phi}(z)$

Monte Carlo
estimate of the
gradient

The Score Function Gradient Estimator

However, there is an issue:

The Score Function Gradient Estimator

However, there is an issue:

- The score-function estimator has an important shortcoming: it has **high variance**.

The Score Function Gradient Estimator

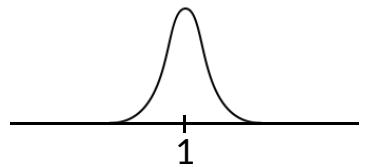
However, there is an issue:

- The score-function estimator has an important shortcoming: it has **high variance**.
- What does this mean?
 - Suppose we are using Monte Carlo to estimate a quantity with expected value = 1.

The Score Function Gradient Estimator

However, there is an issue:

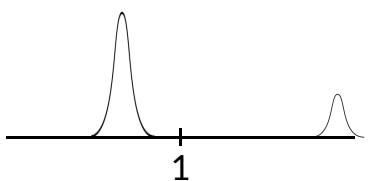
- The score-function estimator has an important shortcoming: it has **high variance**.
- What does this mean?
 - Suppose we are using Monte Carlo to estimate a quantity with expected value = 1.
 - If our samples are: 0.9, 1.1, 0.96, 1.05 (i.e., close to 1), then after a few samples we'll get a good estimate of the true expected value.



The Score Function Gradient Estimator

However, there is an issue:

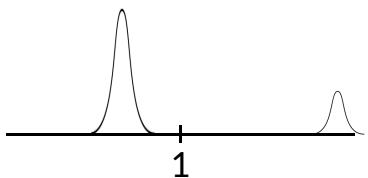
- The score-function estimator has an important shortcoming: it has **high variance**.
- What does this mean?
 - Suppose we are using Monte Carlo to estimate a quantity with expected value = 1.
 - If our samples are: 0.9, 1.1, 0.96, 1.05 (i.e., close to 1), then after a few samples we'll get a good estimate of the true expected value.
 - If our samples are: zero 99% of the time, and one-hundred 1% of the time \Rightarrow expected value is correct, but estimate is worse! (You need many more samples).



The Score Function Gradient Estimator

However, there is an issue:

- The score-function estimator has an important shortcoming: it has **high variance**.
- What does this mean?
 - Suppose we are using Monte Carlo to estimate a quantity with expected value = 1.
 - If our samples are: 0.9, 1.1, 0.96, 1.05 (i.e., close to 1), then after a few samples we'll get a good estimate of the true expected value.
 - If our samples are: zero 99% of the time, and one-hundred 1% of the time \Rightarrow expected value is correct, but estimate is worse! (You need many more samples).
- This latter issue is an example of a high variance MC estimator.



The Score Function Gradient Estimator → SGVB Estimator

One key contribution of the VAE paper is to propose an alternative estimator, which is better behaved (lower variance).

The Score Function Gradient Estimator → SGVB Estimator

One key contribution of the VAE paper is to propose an alternative estimator, which is better behaved (lower variance).

Sometimes referred to as the *stochastic gradient variational Bayes (SGVB) estimator*.

The Score Function Gradient Estimator → SGVB Estimator

One key contribution of the VAE paper is to propose an alternative estimator, which is better behaved (lower variance).

Sometimes referred to as the *stochastic gradient variational Bayes (SGVB) estimator*.

This is done in two steps:

- (1) We reformulate the ELBO so that parts of it can be computed in closed form (*i.e.*, without Monte Carlo).
- (2) Then we use an alternative gradient estimator (SGVB) based on the so-called **reparameterization trick**.

Reformulating the ELBO

The reformulation of the ELBO is as follows:

Reformulating the ELBO

The reformulation of the ELBO is as follows:

$$\log p(x) \geq \mathbb{E}_{q_\phi(z|x)} [\log p_\theta(x | z)] - \text{KL} [q_\phi(z | x) || p_\theta(z)]$$

Reformulating the ELBO

The reformulation of the ELBO is as follows:

$$\log p(x) \geq \mathbb{E}_{q_\phi(z|x)} [\log p_\theta(x | z)] - \text{KL} [q_\phi(z | x) || p_\theta(z)]$$

Evidence

Reformulating the ELBO

The reformulation of the ELBO is as follows:

$$\log p(x) \geq \mathbb{E}_{q_\phi(z|x)} [\log p_\theta(x | z)] - \text{KL} [q_\phi(z | x) || p_\theta(z)]$$

Evidence

ELBO (lower bound on the evidence), consisting of two terms:

Reformulating the ELBO

The reformulation of the ELBO is as follows:

$$\log p(x) \geq \mathbb{E}_{q_\phi(z|x)} [\log p_\theta(x | z)] - \text{KL} [q_\phi(z | x) || p_\theta(z)]$$

Evidence

Expectation of the *log-likelihood*

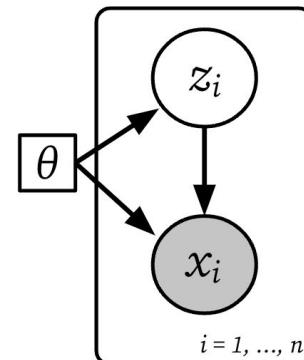
Reformulating the ELBO

The reformulation of the ELBO is as follows:

$$\log p(x) \geq \mathbb{E}_{q_\phi(z|x)} [\log p_\theta(x | z)] - \text{KL} [q_\phi(z | x) || p_\theta(z)]$$

Evidence

Expectation of the *log-likelihood*



We know the log-likelihood (typically chosen when we define our probabilistic model).

Reformulating the ELBO

The reformulation of the ELBO is as follows:

$$\log p(x) \geq \mathbb{E}_{q_\phi(z|x)} [\log p_\theta(x | z)] - \text{KL} [q_\phi(z | x) || p_\theta(z)]$$

Evidence

Expectation of the *log-likelihood*

KL between *q posterior* and *p prior*

Reformulating the ELBO

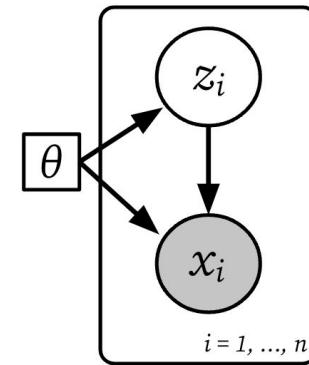
The reformulation of the ELBO is as follows:

$$\log p(x) \geq \mathbb{E}_{q_\phi(z|x)} [\log p_\theta(x | z)] - \text{KL} [q_\phi(z | x) || p_\theta(z)]$$

Evidence

Expectation of the *log-likelihood*

KL between *q posterior* and *p prior*



The prior is also known/chosen when we first define our model.

Reformulating the ELBO

The reformulation of the ELBO is as follows:

$$\log p(x) \geq \mathbb{E}_{q_\phi(z|x)} [\log p_\theta(x | z)] - \text{KL} [q_\phi(z | x) || p_\theta(z)]$$

Evidence

Expectation of the *log-likelihood*

KL between *q posterior* and *p prior*

(Can verify this using some algebra on the previous ELBO formulation).

Reformulating the ELBO

The reformulation of the ELBO is as follows:

$$\log p(x) \geq \mathbb{E}_{q_\phi(z|x)} [\log p_\theta(x | z)] - \text{KL} [q_\phi(z | x) || p_\theta(z)]$$

Evidence

Expectation of the *log-likelihood*

KL between *q posterior* and *p prior*

This has an interesting interpretation!

Reformulating the ELBO

$$\log p(x) \geq \mathbb{E}_{q_\phi(z|x)} [\log p_\theta(x | z)] - \text{KL} [q_\phi(z | x) || p_\theta(z)]$$

Evidence

Expectation of the *log-likelihood*

KL between *q posterior* and *p prior*

Reformulating the ELBO

$$\log p(x) \geq \mathbb{E}_{q_\phi(z|x)} [\log p_\theta(x | z)] - \text{KL} [q_\phi(z | x) || p_\theta(z)]$$

Evidence

Expectation of the *log-likelihood*

KL between *q posterior* and *p prior*

First:

Reformulating the ELBO

$$\log p(x) \geq \mathbb{E}_{q_\phi(z|x)} [\log p_\theta(x | z)] - \text{KL} [q_\phi(z | x) || p_\theta(z)]$$

Evidence

Expectation of the *log-likelihood*

KL between *q posterior* and *p prior*

First:

- Think of x as an observed data point.

Reformulating the ELBO

$$\log p(x) \geq \mathbb{E}_{q_\phi(z|x)} [\log p_\theta(x | z)] - \text{KL} [q_\phi(z | x) || p_\theta(z)]$$

Evidence

Expectation of the *log-likelihood*

KL between *q posterior* and *p prior*

First:

- Think of x as an observed data point.
- Think of $z \sim q_\phi(z | x)$ as a “code” that describes x .
 - (Note that both terms in the right hand side involve an expectation with respect to q .)

Reformulating the ELBO

$$\log p(x) \geq \mathbb{E}_{q_\phi(z|x)} [\log p_\theta(x | z)] - \text{KL} [q_\phi(z | x) || p_\theta(z)]$$

Evidence

Expectation of the *log-likelihood*

KL between *q posterior* and *p prior*

First:

- Think of x as an observed data point.
- Think of $z \sim q_\phi(z | x)$ as a “code” that describes x .
 - (Note that both terms in the right hand side involve an expectation with respect to q .)
- ⇒ We will therefore call q_ϕ the **encoder**.

Reformulating the ELBO

$$\log p(x) \geq \mathbb{E}_{q_\phi(z|x)} [\log p_\theta(x | z)] - \text{KL} [q_\phi(z | x) || p_\theta(z)]$$

Evidence

Expectation of the *log-likelihood*

KL between *q posterior* and *p prior*

Next:

Reformulating the ELBO

$$\log p(x) \geq \mathbb{E}_{q_\phi(z|x)} [\log p_\theta(x | z)] - \text{KL} [q_\phi(z | x) || p_\theta(z)]$$

Evidence

Expectation of the *log-likelihood*

KL between *q posterior* and *p prior*

Next:

encoder

- The first term is the log-likelihood of an x , given the code $z \sim q_\phi(z | x)$ conditioned on x .

Reformulating the ELBO

$$\log p(x) \geq \mathbb{E}_{q_\phi(z|x)} [\log p_\theta(x | z)] - \text{KL} [q_\phi(z | x) || p_\theta(z)]$$

Evidence

Expectation of the *log-likelihood*

KL between *q posterior* and *p prior*

Next:

encoder

- The first term is the log-likelihood of an x , given the code $z \sim q_\phi(z | x)$ conditioned on x .
- ⇒ this first term is maximized when $p_\theta(x | z)$ assigns high probability to this x .

Reformulating the ELBO

$$\log p(x) \geq \mathbb{E}_{q_\phi(z|x)} [\log p_\theta(x | z)] - \text{KL} [q_\phi(z | x) || p_\theta(z)]$$

Evidence

Expectation of the *log-likelihood*

KL between *q posterior* and *p prior*

Next:

encoder

- The first term is the log-likelihood of an x , given the code $z \sim q_\phi(z | x)$ conditioned on x .
- \Rightarrow this first term is maximized when $p_\theta(x | z)$ assigns high probability to this x .
- \Rightarrow i.e., It is trying to reconstruct x given the code z (drawn from $q_\phi(z | x)$).

Reformulating the ELBO

$$\log p(x) \geq \mathbb{E}_{q_\phi(z|x)} [\log p_\theta(x | z)] - \text{KL} [q_\phi(z | x) || p_\theta(z)]$$

Evidence

Expectation of the *log-likelihood*

KL between *q posterior* and *p prior*

Next:

- The first term is the log-likelihood of an x , given the code $z \sim q_\phi(z | x)$ conditioned on x .
- \Rightarrow this first term is maximized when $p_\theta(x | z)$ assigns high probability to this x .
- \Rightarrow i.e., It is trying to reconstruct x given the code z (drawn from $q_\phi(z | x)$).
- \Rightarrow We will therefore call $p_\theta(x | z)$ the **decoder**.
 - And the first term we call the *reconstruction error*.

Reformulating the ELBO

$$\log p(x) \geq \mathbb{E}_{q_\phi(z|x)} [\log p_\theta(x | z)] - \text{KL} [q_\phi(z | x) || p_\theta(z)]$$

Evidence

Expectation of the *log-likelihood*

KL between *q posterior* and *p prior*

Also:

Reformulating the ELBO

$$\log p(x) \geq \mathbb{E}_{q_\phi(z|x)} [\log p_\theta(x | z)] - \text{KL} [q_\phi(z | x) || p_\theta(z)]$$

Evidence

Expectation of the *log-likelihood*

KL between *q posterior* and *p prior*

Also:

- The second term is the KL divergence between $q_\phi(z | x)$ and the model's prior $p_\theta(z)$.

Reformulating the ELBO

$$\log p(x) \geq \mathbb{E}_{q_\phi(z|x)} [\log p_\theta(x | z)] - \text{KL} [q_\phi(z | x) || p_\theta(z)]$$

Evidence

Expectation of the *log-likelihood*

KL between *q posterior* and *p prior*

Also:

- The second term is the KL divergence between $q_\phi(z | x)$ and the model's prior $p_\theta(z)$.
- (In this model, we will assume the prior on z is a unit Normal distribution). $\mathcal{N}(0, I)$
 - \Rightarrow this encourages the codes z to look Gaussian distributed.

Reformulating the ELBO

$$\log p(x) \geq \mathbb{E}_{q_\phi(z|x)} [\log p_\theta(x | z)] - \text{KL} [q_\phi(z | x) || p_\theta(z)]$$

Evidence

Expectation of the *log-likelihood*

KL between *q posterior* and *p prior*

Also:

- The second term is the KL divergence between $q_\phi(z | x)$ and the model's prior $p_\theta(z)$.
- (In this model, we will assume the prior on z is a unit Normal distribution). $\mathcal{N}(0, I)$
 - \Rightarrow this encourages the codes z to look Gaussian distributed.
- We call this second term the *regularization term*.
 - It prevents $q_\phi(z | x)$ from, e.g., simply encoding an identity mapping.

Reformulating the ELBO

$$\log p(x) \geq \mathbb{E}_{q_\phi(z|x)} [\log p_\theta(x | z)] - \text{KL} [q_\phi(z | x) || p_\theta(z)]$$

Evidence

Expectation of the *log-likelihood*

KL between *q posterior* and *p prior*

In summary:

Reformulating the ELBO

$$\log p(x) \geq \mathbb{E}_{q_\phi(z|x)} [\log p_\theta(x | z)] - \text{KL} [q_\phi(z | x) || p_\theta(z)]$$

Evidence

Expectation of the *log-likelihood*

KL between *q posterior* and *p prior*

In summary:

When we optimize this ELBO, we are trying to fit a $q_\phi(z | x)$ that will map an observation x onto a useful latent variable z ...

... and also to fit a $p_\theta(x | z)$ which reconstructs this x , given the latent z .

Reformulating the ELBO

$$\log p(x) \geq \mathbb{E}_{q_\phi(z|x)} [\log p_\theta(x | z)] - \text{KL} [q_\phi(z | x) || p_\theta(z)]$$

Evidence

Expectation of the *log-likelihood*

KL between *q posterior* and *p prior*

In summary:

When we optimize this ELBO, we are trying to fit a $q_\phi(z | x)$ that will map an observation x onto a useful latent variable z ...

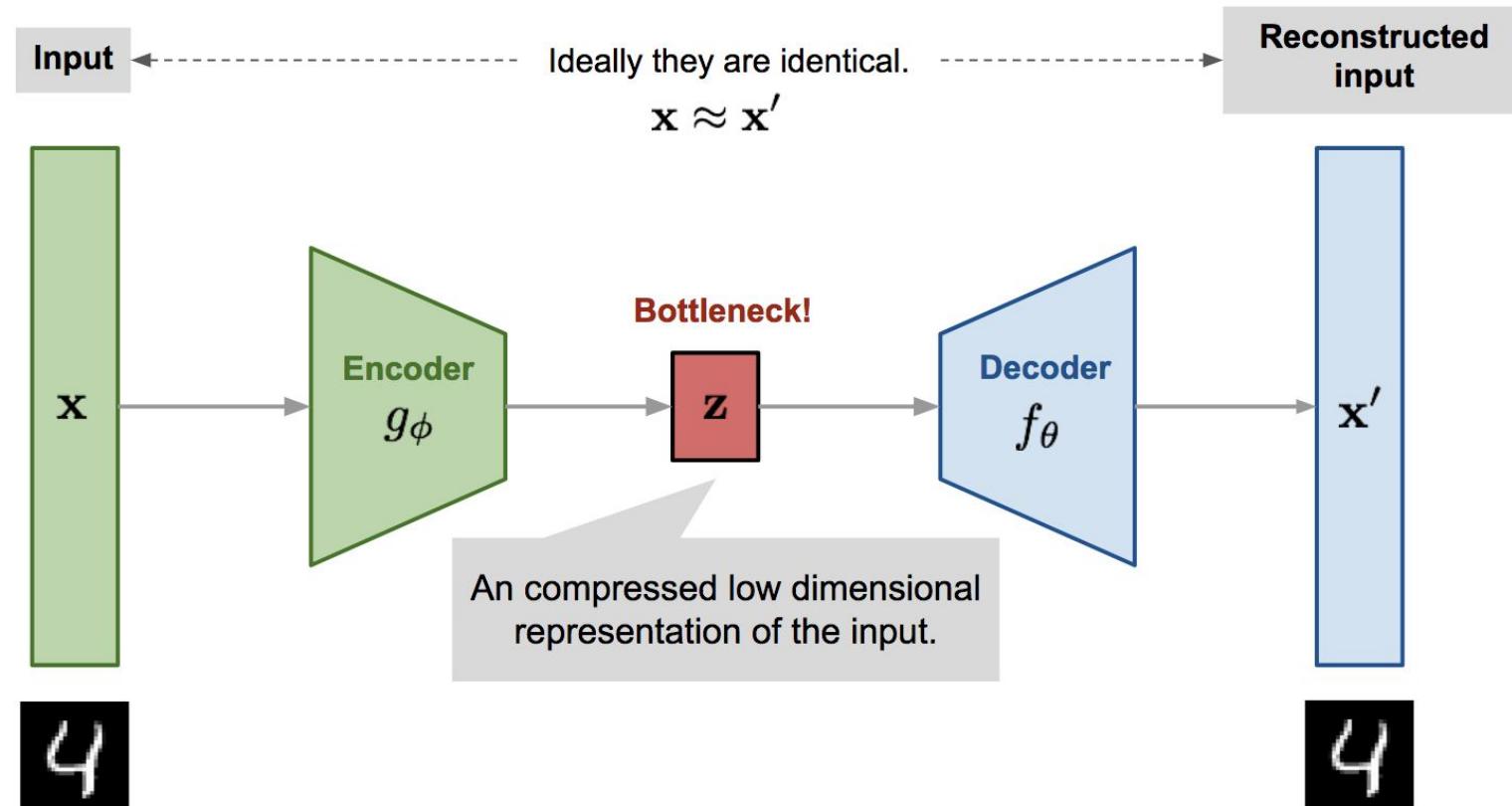
... and also to fit a $p_\theta(x | z)$ which reconstructs this x , given the latent z .

⇒ This is reminiscent of an *autoencoder neural network*.

(i.e., a pair of neural networks, which encode/decode data, aiming to minimize the reconstruction error).

Reconstruction Term \Leftrightarrow Autoencoder

Aside – illustration of an **autoencoder** (first developed in 1990s, or even late 1980s).



Reparameterization Trick

So how can we use this reformulation to form a *lower-variance gradient estimator*, and then optimize the ELBO?

Reparameterization Trick

So how can we use this reformulation to form a *lower-variance gradient estimator*, and then optimize the ELBO?

We will use something call the **reparameterization trick**!

Reparameterization Trick

So how can we use this reformulation to form a *lower-variance gradient estimator*, and then optimize the ELBO?

We will use something call the **reparameterization trick**!

Recall:

This is done in two steps:

SGVB Estimator

- (1) We reformulate the ELBO  .
- (2) Then we use an alternative gradient estimator (SGVB) based on the so-called **reparameterization trick**.

Reparameterization Trick

Under certain mild conditions, we can express $q_\phi(z \mid x)$ via the following two-step generative process:

Reparameterization Trick

Under certain mild conditions, we can express $q_\phi(z \mid x)$ via the following two-step generative process:

- First sample a noise variable ϵ from a simple distribution $p(\epsilon)$ (such as a standard Normal $\mathcal{N}(0, I)$):

$$\epsilon \sim p(\epsilon)$$

Note: this does not depend on ϕ .

Reparameterization Trick

Under certain mild conditions, we can express $q_\phi(z \mid x)$ via the following two-step generative process:

- First sample a noise variable ϵ from a simple distribution $p(\epsilon)$ (such as a standard Normal $\mathcal{N}(0, I)$):

$$\epsilon \sim p(\epsilon)$$

Note: this does not depend on ϕ .

- Then we apply a deterministic transformation $g_\phi(\epsilon, x)$ that maps the random noise into a more-complex distribution:

$$z = g_\phi(\epsilon, x)$$

Reparameterization Trick

Note: for many interesting classes of q_ϕ , it will be possible to define a g_ϕ such that $z = g_\phi(\epsilon, x)$ is distributed according to $q_\phi(z \mid x)$.

Reparameterization Trick

Note: for many interesting classes of q_ϕ , it will be possible to define a g_ϕ such that $z = g_\phi(\epsilon, x)$ is distributed according to $q_\phi(z \mid x)$.

Graphical model of $q_\phi(z \mid x)$
can be viewed as similar to...
(recall from first PGM lecture)

Famous Bayesian Networks – Neural Transformations in PGMs

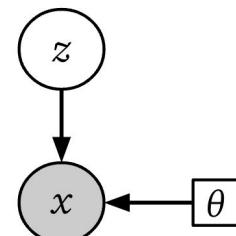
Can incorporate more-complex function transformations (or neural networks) into Bayesian networks!

But instead we could do:

$$z \sim p(z) = \mathcal{N}(0, 1)$$

$$x \sim p_\theta(x \mid z) = \mathcal{N}(g_\theta(z), 1)$$

$g_\theta(z)$ might be a complex fixed/known, or have parameters that we learn (e.g., if we use a neural network!).



Reparameterization Trick

We refer to this as a ***reparameterization trick*** (we are reparameterizing, i.e., writing out $q_\phi(z \mid x)$ in a different way).

Reparameterization Trick

We refer to this as a **reparameterization trick** (we are reparameterizing, i.e., writing out $q_\phi(z | x)$ in a different way).

Simplest example of the reparameterization trick:

Reparameterization Trick

We refer to this as a **reparameterization trick** (we are reparameterizing, i.e., writing out $q_\phi(z | x)$ in a different way).

Simplest example of the reparameterization trick:

- Instead of writing: $z \sim q_{\mu, \sigma}(z) = \mathcal{N}(z; \mu, \sigma^2)$

Parameter ϕ here is: (μ, σ) .

Reparameterization Trick

We refer to this as a **reparameterization trick** (we are reparameterizing, i.e., writing out $q_\phi(z | x)$ in a different way).

Simplest example of the reparameterization trick:

Parameter ϕ here is: (μ, σ) .

- Instead of writing: $z \sim q_{\mu, \sigma}(z) = \mathcal{N}(z; \mu, \sigma^2)$
- We can write: $z = g_{\mu, \sigma}(\epsilon) = \mu + \epsilon \cdot \sigma$, where $\epsilon \sim \mathcal{N}(\epsilon ; 0, 1)$

Reparameterization Trick

Main advantages of the reparameterization trick:

Reparameterization Trick

Main advantages of the reparameterization trick:

We can now write the gradient of an expectation (gradient wrt ϕ , expectation of q), for a given test function f , as:

Reparameterization Trick

Main advantages of the reparameterization trick:

We can now write the gradient of an expectation (gradient wrt ϕ , expectation of q), for a given test function f , as:

$$\begin{aligned}\nabla_{\phi} \mathbb{E}_{z \sim q_{\phi}(z|x)} [f(x, z)] &= \nabla_{\phi} \mathbb{E}_{\epsilon \sim p(\epsilon)} [f(x, g_{\phi}(\epsilon, x))] \\ &= \mathbb{E}_{\epsilon \sim p(\epsilon)} [\nabla_{\phi} f(x, g_{\phi}(\epsilon, x))]\end{aligned}$$

Reparameterization Trick

Main advantages of the reparameterization trick:

We can now write the gradient of an expectation (gradient wrt ϕ , expectation of q), for a given test function f , as:

$$\begin{aligned}\nabla_{\phi} \mathbb{E}_{z \sim q_{\phi}(z|x)} [f(x, z)] &= \nabla_{\phi} \mathbb{E}_{\epsilon \sim p(\epsilon)} [f(x, g_{\phi}(\epsilon, x))] \\ &= \mathbb{E}_{\epsilon \sim p(\epsilon)} [\nabla_{\phi} f(x, g_{\phi}(\epsilon, x))]\end{aligned}$$

Replace q with $p(\epsilon)$,
and z with $g(\epsilon, x)$.

Reparameterization Trick

Main advantages of the reparameterization trick:

We can now write the gradient of an expectation (gradient wrt ϕ , expectation of q), for a given test function f , as:

$$\nabla_{\phi} \mathbb{E}_{z \sim q_{\phi}(z|x)} [f(x, z)] = \nabla_{\phi} \mathbb{E}_{\epsilon \sim p(\epsilon)} [f(x, g_{\phi}(\epsilon, x))]$$

$$= \mathbb{E}_{\epsilon \sim p(\epsilon)} [\nabla_{\phi} f(x, g_{\phi}(\epsilon, x))]$$

Replace q with $p(\epsilon)$,
and z with $g(\epsilon, x)$.

Can now swap
gradient and
expectation!

Reparameterization Trick

Main advantages of the reparameterization trick:

We can now write the gradient of an expectation (gradient wrt ϕ , expectation of q), for a given test function f , as:

$$\begin{aligned}\nabla_{\phi} \mathbb{E}_{z \sim q_{\phi}(z|x)} [f(x, z)] &= \nabla_{\phi} \mathbb{E}_{\epsilon \sim p(\epsilon)} [f(x, g_{\phi}(\epsilon, x))] \\ &= \mathbb{E}_{\epsilon \sim p(\epsilon)} [\nabla_{\phi} f(x, g_{\phi}(\epsilon, x))]\end{aligned}$$

Replace q with $p(\epsilon)$,
and z with $g(\epsilon, x)$.

Can now swap
gradient and
expectation!

Can now take a Monte Carlo estimate of this expectation.

⇒ yields an estimate of the gradient for ϕ , with much lower variance than the score-function estimator.

What to choose for p and q ?

What to choose for p and q ?

Suppose, after learning θ , we hope to use $p_\theta(x \mid z)$ as a generative model.

(And we hope to use $q_\phi(z \mid x)$ to perform inference in order to learn θ).

What to choose for p and q ?

Suppose, after learning θ , we hope to use $p_\theta(x \mid z)$ as a generative model.

(And we hope to use $q_\phi(z \mid x)$ to perform inference in order to learn θ).

Question: how should we specify the form of p_θ and q_ϕ ?

What to choose for p and q ?

Suppose, after learning θ , we hope to use $p_\theta(x \mid z)$ as a generative model.

(And we hope to use $q_\phi(z \mid x)$ to perform inference in order to learn θ).

Question: how should we specify the form of p_θ and q_ϕ ?

Note:

What to choose for p and q ?

Suppose, after learning θ , we hope to use $p_\theta(x \mid z)$ as a generative model.

(And we hope to use $q_\phi(z \mid x)$ to perform inference in order to learn θ).

Question: how should we specify the form of p_θ and q_ϕ ?

Note:

- We want $p_\theta(x \mid z)$ to be flexible enough to represent the richness of complex/high-dimensional data.

What to choose for p and q ?

Suppose, after learning θ , we hope to use $p_\theta(x \mid z)$ as a generative model.

(And we hope to use $q_\phi(z \mid x)$ to perform inference in order to learn θ).

Question: how should we specify the form of p_θ and q_ϕ ?

Note:

- We want $p_\theta(x \mid z)$ to be flexible enough to represent the richness of complex/high-dimensional data.
- We want $q_\phi(z \mid x)$ to be flexible enough to approximate $p_\theta(z \mid x)$ (i.e., the true posterior) well.

What to choose for p and q ?

Answer: in VAEs, we will parameterize p_θ and q_ϕ using neural networks.



What to choose for p and q ?

Answer: in VAEs, we will parameterize p_θ and q_ϕ using neural networks.



As an example, we can define $q_\phi(z \mid x)$ to be:

$$q_\phi(z \mid x) = \mathcal{N}(z ; \mu_\phi(x), \text{diag}(\sigma_\phi(x))^2)$$

Where $\mu_\phi(x)$ and $\sigma_\phi(x)$ are vector-valued functions of x parametrized by an arbitrary neural network.

What to choose for p and q ?

Answer: in VAEs, we will parameterize p_θ and q_ϕ using neural networks.



As an example, we can define $q_\phi(z \mid x)$ to be:

$$q_\phi(z \mid x) = \mathcal{N}(z ; \mu_\phi(x), \text{diag}(\sigma_\phi(x))^2)$$

Where $\mu_\phi(x)$ and $\sigma_\phi(x)$ are vector-valued functions of x parametrized by an arbitrary neural network.

⇒ could call them: “*mean network*” and “*standard deviation network*”.

What to choose for p and q ?

Answer: in VAEs, we will parameterize p_θ and q_ϕ using neural networks.



As an example, we can define $q_\phi(z \mid x)$ to be:

$$q_\phi(z \mid x) = \mathcal{N}(z ; \mu_\phi(x), \text{diag}(\sigma_\phi(x))^2)$$

More generally, can be applied to any exponential family distribution by parameterizing the sufficient statistics by a function of x .

Where $\mu_\phi(x)$ and $\sigma_\phi(x)$ are vector-valued functions of x parametrized by an arbitrary neural network.

⇒ could call them: “*mean network*” and “*standard deviation network*”.

Variational Autoencoder (VAE)

So what is the variational autoencoder (VAE), and how does it relate to the AEVB algorithm?

Variational Autoencoder (VAE)

First, to summarize – the AEVB algorithm is simply the combination of:

Variational Autoencoder (VAE)

First, to summarize – the AEVB algorithm is simply the combination of:

- (1) The alternative ELBO reformulation.

$$\mathbb{E}_{q_\phi(z|x)} [\log p_\theta(x | z)] - \text{KL} [q_\phi(z | x) || p_\theta(z)]$$

Expectation of the *log-likelihood*

KL between *q posterior* and *p prior*

Variational Autoencoder (VAE)

First, to summarize – the AEVB algorithm is simply the combination of:

- (1) The alternative ELBO reformulation.

$$\mathbb{E}_{q_\phi(z|x)} [\log p_\theta(x | z)] - \text{KL} [q_\phi(z | x) || p_\theta(z)]$$

Expectation of the *log-likelihood*

KL between *q posterior* and *p prior*

- (2) The black-box (gradient-based) variational inference approach.

Variational Autoencoder (VAE)

First, to summarize – the AEVB algorithm is simply the combination of:

- (1) The alternative ELBO reformulation.

$$\mathbb{E}_{q_\phi(z|x)} [\log p_\theta(x | z)] - \text{KL} [q_\phi(z | x) || p_\theta(z)]$$

Expectation of the *log-likelihood*

KL between *q posterior* and *p prior*

- (2) The black-box (gradient-based) variational inference approach.
- (3) The reparameterization-based SGVB (low-variance) estimator.

Variational Autoencoder (VAE)

First, to summarize – the AEVB algorithm is simply the combination of:

- (1) The alternative ELBO reformulation.

$$\mathbb{E}_{q_\phi(z|x)} [\log p_\theta(x | z)] - \text{KL} [q_\phi(z | x) || p_\theta(z)]$$

Expectation of the *log-likelihood*

KL between *q posterior* and *p prior*

- (2) The black-box (gradient-based) variational inference approach.
- (3) The reparameterization-based SGVB (low-variance) estimator.

⇒ takes gradient steps on θ and ϕ to optimize the alternative ELBO.

Variational Autoencoder (VAE)

The VAE is a special case – using the AEVB algorithm for a specific choice of p and q .

Variational Autoencoder (VAE)

The **VAE** is a special case – using the AEVB algorithm for a specific choice of p and q .

The model p_{θ} is parameterized as:

Variational Autoencoder (VAE)

The VAE is a special case – using the AEVB algorithm for a specific choice of p and q .

The model p_θ is parameterized as:

$$p_\theta(z) = \mathcal{N}(z; 0, I)$$
$$p_\theta(x | z) = \mathcal{N}(x; \mu_\theta(z), \text{diag}(\sigma_\theta(z))^2)$$

Variational Autoencoder (VAE)

The VAE is a special case – using the AEVB algorithm for a specific choice of p and q .

The model p_θ is parameterized as:

$$p_\theta(z) = \mathcal{N}(z; 0, I)$$

$$p_\theta(x | z) = \mathcal{N}(x; \mu_\theta(z), \text{diag}(\sigma_\theta(z))^2)$$

Prior on latent variable
 z is a Gaussian.

Variational Autoencoder (VAE)

The VAE is a special case – using the AEVB algorithm for a specific choice of p and q .

The model p_θ is parameterized as:

$$p_\theta(z) = \mathcal{N}(z; 0, I)$$

$$p_\theta(x | z) = \mathcal{N}(x; \mu_\theta(z), \text{diag}(\sigma_\theta(z))^2)$$

Prior on latent variable
 z is a Gaussian.

Decode latent variable z
into data point x .
(Also Gaussian).

Variational Autoencoder (VAE)

The **VAE** is a special case – using the AEVB algorithm for a specific choice of p and q .

The posterior approximation q_ϕ is parameterized as:

Variational Autoencoder (VAE)

The VAE is a special case – using the AEVB algorithm for a specific choice of p and q .

The posterior approximation q_ϕ is parameterized as:

$$q_\phi(z | x) = \mathcal{N}(z; \mu_\phi(x), \text{diag}(\sigma_\phi(x))^2)$$

Encode data point x as
latent variable z .
(Also Gaussian).

Variational Autoencoder (VAE)

Useful properties of this choice of p_θ and q_ϕ :

Variational Autoencoder (VAE)

Useful properties of this choice of p_θ and q_ϕ :

- Closed-form expression to compute the regularization term (2nd term)

$$\mathbb{E}_{q_\phi(z|x)} [\log p_\theta(x | z)] - \text{KL} [q_\phi(z | x) || p_\theta(z)]$$

Expectation of the *log-likelihood*

KL between *q posterior* and *p prior*

- KL divergence of two Gaussians \Rightarrow closed form expression and gradient!

Variational Autoencoder (VAE)

Useful properties of this choice of p_θ and q_ϕ :

- Closed-form expression to compute the regularization term (2nd term)

$$\mathbb{E}_{q_\phi(z|x)} [\log p_\theta(x | z)] - \text{KL} [q_\phi(z | x) || p_\theta(z)]$$

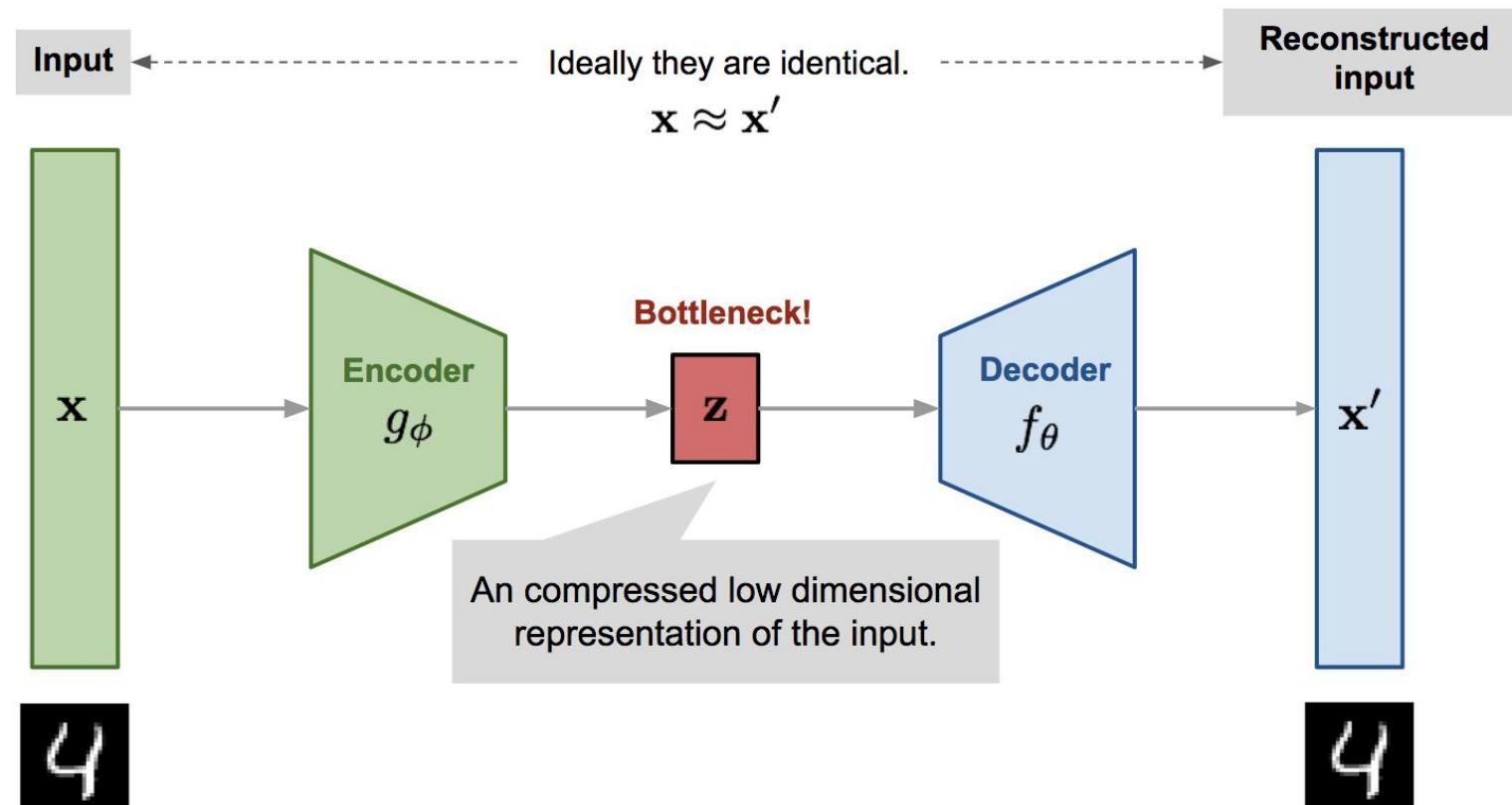
Expectation of the *log-likelihood*KL between *q posterior* and *p prior*



- KL divergence of two Gaussians \Rightarrow closed form expression and gradient!
- So we only need to perform Monte Carlo gradient step on the 1st term.
 - And can carry out SGVB estimator using neural network reparameterization trick as described above.

Variational Autoencoder (VAE)

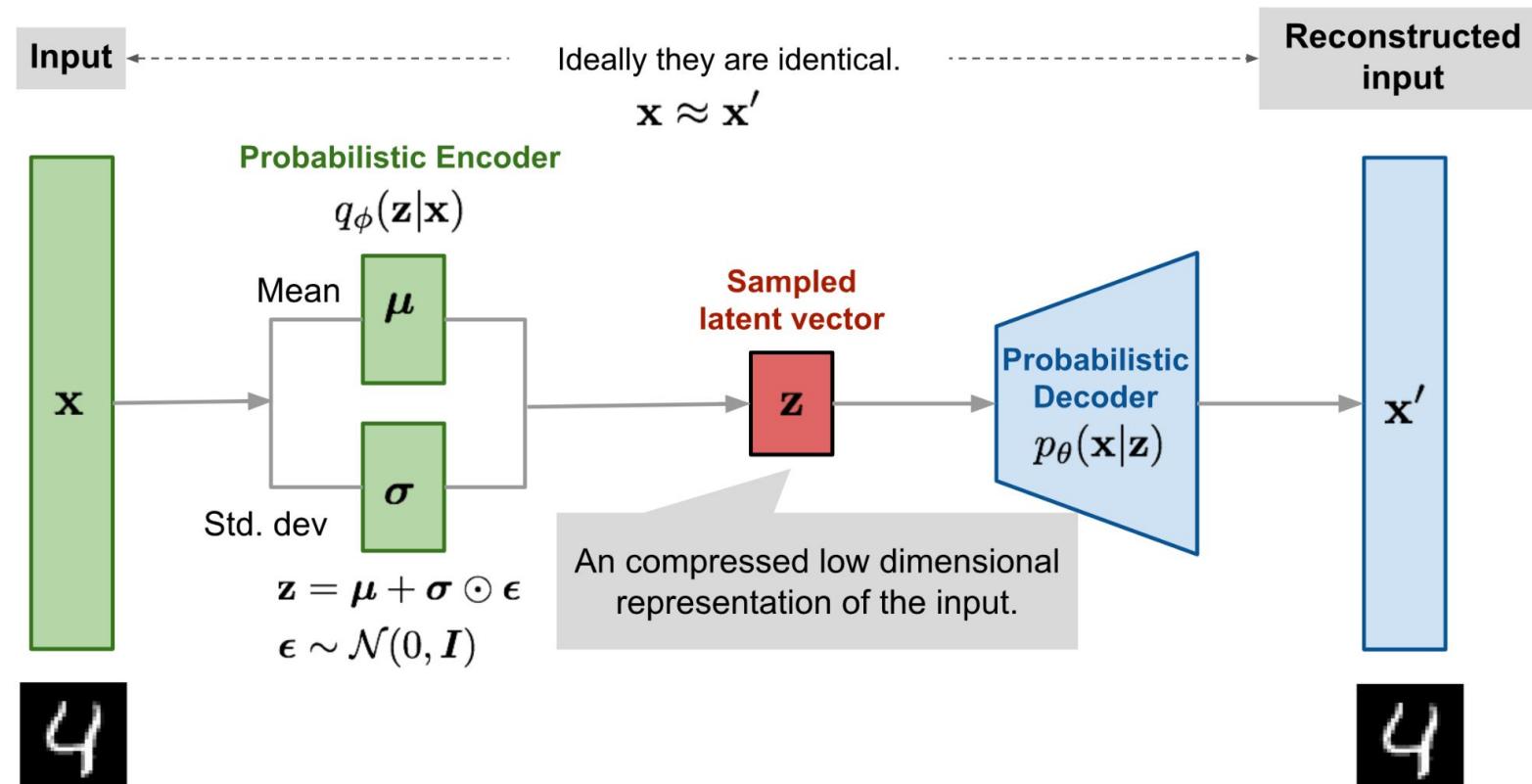
Visual summary – back to the *autoencoder*:



Source: Lil'Log,
"From Autoencoder
to Beta-VAE"

Variational Autoencoder (VAE)

Visual summary – instead, a *variational autoencoder*:



Variational Autoencoder (VAE) – Results

So does it work?

Variational Autoencoder (VAE) – Results

So does it work? → Results from the original paper (*ICLR 2014*)

Variational Autoencoder (VAE) – Results

So does it work? → Results from the original paper (*ICLR 2014*)

Auto-Encoding Variational Bayes

Diederik P. Kingma
Machine Learning Group
Universiteit van Amsterdam
dpkingma@gmail.com

Max Welling
Machine Learning Group
Universiteit van Amsterdam
welling.max@gmail.com

Abstract

How can we perform efficient inference and learning in directed probabilistic models, in the presence of continuous latent variables with intractable posterior distributions? The variational Bayesian (VB) approach involves the optimization of an approximation to the intractable posterior. Unfortunately, the common mean-field approach requires analytical solutions of expectations w.r.t. the approximate posterior, which are also intractable in the general case. We show how a reparameterization of the variational lower bound yields a simple differentiable unbiased estimator of the lower bound; this SGVB (Stochastic Gradient Variational Bayes) estimator can be used for efficient approximate posterior inference in almost any model with continuous latent variables and/or parameters, and is straightforward to optimize using standard stochastic gradient ascent techniques.

For the case of an i.i.d. dataset and a linear latent variable per datapoint, we propose the Auto-Encoding VB (AEVB) algorithm. In the AEVB algorithm we perform inference and learning simultaneously by using the SGVB estimator to optimize a recognition model that allows us to perform very efficient approximate posterior inference using simple ancestral sampling, which in turn allows us to efficiently learn the model parameters, without the need of expensive iterative inference schemes (such as MCMC) per datapoint. The learned approximate posterior inference model can also be used for a host of tasks such as recognition, denoising, representation and visualization purposes. When a neural network is used for the recognition model, we arrive at the *variational auto-encoder*.

1 Introduction

How can we perform efficient approximate inference and learning with directed probabilistic models whose continuous latent variables and/or parameters have intractable posterior distributions? The variational Bayesian (VB) approach involves the optimization of an approximation to the intractable posterior. Unfortunately, the common mean-field approach requires analytical solutions of expectations w.r.t. the approximate posterior, which are also intractable in the general case. We show how a reparameterization of the variational lower bound yields a simple differentiable unbiased estimator of the lower bound; this SGVB (Stochastic Gradient Variational Bayes) estimator can be used for efficient approximate posterior inference in almost any model with continuous latent variables and/or parameters, and is straightforward to optimize using standard stochastic gradient ascent techniques.

For the case of an i.i.d. dataset and a linear latent variable per datapoint, we propose the Auto-Encoding VB (AEVB) algorithm. In the AEVB algorithm we perform inference and learning simultaneously by using the SGVB estimator to optimize a recognition model that allows us to perform very efficient approximate posterior inference using simple ancestral sampling, which in turn allows us to efficiently learn the model parameters, without the need of expensive iterative inference schemes (such as MCMC) per datapoint. The learned approximate posterior inference model can also be used for a host of tasks such as recognition, denoising, representation and visualization purposes. When a neural network is used for the recognition model, we arrive at the *variational auto-encoder*.

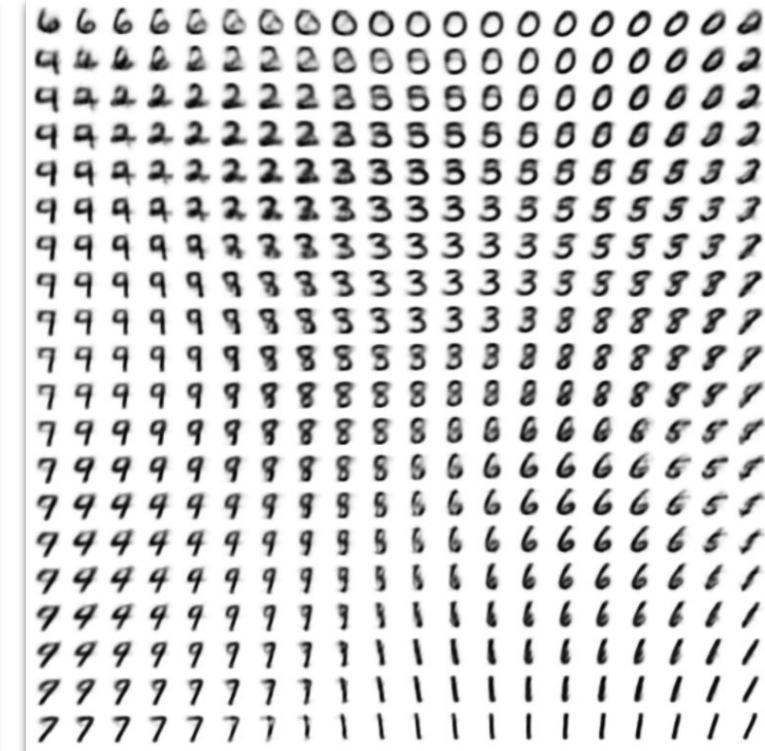
2 Method

The strategy in this section can be used to derive a lower bound estimator (a stochastic objective function) for a variety of directed graphical models with continuous latent variables. We will restrict ourselves here to the common case where we have an i.i.d. dataset with latent variables, a dataset, and where we like to perform maximum likelihood (ML) or maximum a posteriori (MAP) inference on the (global) parameters, and variational inference on the latent variables. It is, for example,

1



Frey Faces Dataset



MNIST Dataset

Variational Autoencoder (VAE) – Results

So does it work? → A few years later ... (*NeurIPS 2019*)

Variational Autoencoder (VAE) – Results

So does it work? → A few years later ... (NeurIPS 2019)

Generating Diverse High-Fidelity Images with VQ-VAE-2

Ali Razavi^{*} DeepMind alirazavi@google.com Aaron van den Oord^{*} DeepMind avdnoord@google.com Oriol Vinyals DeepMind vinyals@google.com

Abstract

We explore the use of Vector Quantized Variational AutoEncoder (VQ-VAE) models for large scale image generation. To this end, we scale and enhance the autoregressive prior used in VQ-VAE to generate a heterogenous set of images hierarchically and at multiple scales, faster than previous methods. We use simple feed-forward encoder and decoder networks, making our model an attractive candidate for applications where the encoding and/or decoding speed is critical. Additionally, VQ-VAE requires sampling an autoregressive model only in the compressed latent space, which is much smaller than the full image in the pixel space, especially for large images. We demonstrate that a multi-scale hierarchical generation of VQ-VAE, augmented with powerful priors over the latent codes, is able to generate samples with quality that rivals that of state of the art Generative Adversarial Networks on multi-faceted datasets such as ImageNet, while not suffering from GAN's known shortcomings such as mode collapse and lack of diversity.

1 Introduction

Deep generative models have significantly improved in the past few years [5, 27, 25]. This is, in part, thanks to architectural innovations as well as computation advances that allows training them at larger scale in both amount of data and model size. The samples generated from these models are hard to distinguish from real data without close inspection, and their applications range from super resolution [21] to domain editing [44], artistic manipulation [36], or text-to-speech and music generation [25].



Figure 1: Class-conditional 256x256 image samples from a two-level model trained on ImageNet.

We distinguish two main types of generative models: likelihood based models, which include VAEs [16, 31], flow based [9, 30, 10, 17] and autoregressive models [20, 39]; and implicit generative

*Equal contributions.

33rd Conference on Neural Information Processing Systems (NeurIPS 2019), Vancouver, Canada.



FFHQ Dataset

Source: Razavi, Ali, Aaron Van den Oord, and Oriol Vinyals. "Generating diverse high-fidelity images with vq-vae-2", 2019.

Intro to Score-based Generative Models

VAEs → Score-based/Diffusion Generative Models

VAEs → Score-based/Diffusion Generative Models

VAEs and GANs were the two most-popular generative models for a few years, until around 2019*.

When a new paradigm of generative modeling came around...

**Actually, the first diffusion paper was in 2015 (or before), but it didn't really gain much attention until around 2019.*

Score-based/Diffusion Generative Models

Score-based generative models, and diffusion generative models:

Deep Unsupervised Learning using Nonequilibrium Thermodynamics

Jascha Sohl-Dickstein
Stanford University

Eric A. Weiss
University of California, Berkeley

Niru Maheswaranathan
Stanford University

Surya Ganguli
Stanford University

Abstract

A central problem in machine learning involves modeling complex data-sets using highly flexible models of the underlying distribution. Data, training, sampling, inference, and evaluation are still analytically or computationally tractable. Here, we develop an approach that simultaneously achieves both flexibility and tractability. The essential idea, inspired by non-equilibrium statistical mechanics, is to learn and then destroy structure in a data distribution through an iterative forward diffusion process. We then learn a reverse diffusion process that restores structure in data, yielding a highly flexible and tractable generative model of the data. This approach allows us to train, sample, infer, and evaluate probabilities in deep generative models with thousands of layers or time steps, as well as to compute conditional and posterior probabilities under the learned model. We additionally release an open source reference implementation of the algorithm.

1. Introduction

Historically, probabilistic models suffer from a tradeoff between two conflicting objectives: *tractability* and *flexibility*. Models that are *tractable* can be analytically evaluated and easily fit to data (e.g. a Gaussian or Laplace). However, non-parametric methods (Gershman & Blei, 2012) can also be very effective¹.

1.1. Diffusion probabilistic models

We present a novel way to define probabilistic models that allows:

1. extreme flexibility in model structure,
2. exact sampling.

¹Non-parametric methods can be seen as transitioning smoothly between tractable and flexible models. For instance, a non-parametric Gaussian mixture model will represent a small amount of data using a single Gaussian, but may represent infinite data as a mixture of an infinite number of Gaussians.

arXiv:1503.03585v8 [cs.LG] 18 Nov 2015

Generative Modeling by Estimating Gradients of the Data Distribution

Yang Song
Stanford University
yangsong@cs.stanford.edu

Stefano Ermon
Stanford University
ermon@cs.stanford.edu

Abstract

We introduce a new generative model where samples are produced via Langevin dynamics using gradients of the data distribution estimated with score matching. Because gradients can be ill-defined and hard to estimate when the data resides on low-dimensional manifolds, we perform the data denoising with two parallel Gaussian noise, and jointly estimate the corresponding scores (i.e., the vector fields of gradients of the perturbed data distribution for all noise levels). For sampling, we propose an annealed Langevin dynamics where we use gradients corresponding to gradually increasing noise levels in the score matching process to close the loop. One major advantage of this flexible model architecture, relative to sampling during training or the use of adversarial methods, is that it provides a learning objective that can be used for principled model comparisons. Our models produce samples comparable to GANs on MNIST, CelebA and CIFAR-10 datasets, achieving a new state-of-the-art inception score of 8.87 on CIFAR-10. Additionally, we demonstrate that our models learn effective representations via image inpainting experiments.

1 Introduction

Generative models have many applications in machine learning. To list a few, they have been used to generate high-fidelity images [26, 6], synthesize realistic speech and music fragments [58], improve the performance of semi-supervised learning [20], generate realistic samples and other anomalous data [54], image learning [22], and explore promising states in reinforcement learning [41]. Recent progress is mainly driven by two approaches: likelihood-based methods [17, 29, 11, 60] and generative adversarial networks (GAN [15]). The former uses log-likelihood (or a suitable surrogate) as the training objective, while the latter uses adversarial training to minimize $J(\theta, \phi)$ [15] (see [20] for a survey).

Although likelihood-based models and GANs have achieved great success, they have some intrinsic limitations. For example, likelihood-based models either have to use specialized architectures to build a normalized probability model (e.g., autoregressive models, flow models), or use surrogate losses (e.g., the evidence lower bound used in variational auto-encoders [29]), contrastive divergence in energy-based models [21] for training. GANs are some of the illustrations of likelihood-based models, but they are not suitable for evaluating and comparing different GAN models. While other objectives exist for generative modeling, such as noise contrastive estimation [19] and minimum probability flow [50], these methods typically only work well for low-dimensional data.

In this paper, we explore a new paradigm for generative modeling based on estimating gradients from the log-density [33] of the logarithmic data density, which is the gradient of the log-density function at the input data point. This is a vector field pointing in the direction where the log data density grows the most. We use a neural network trained with score matching [24] to learn this vector field from data. We then produce samples using Langevin dynamics, which approximately

arXiv:1907.05600v3 [cs.LG] 10 Oct 2020

Denoising Diffusion Probabilistic Models

Jonathan Ho
UC Berkeley
jonathanho@berkeley.edu

Ajay Jain
UC Berkeley
ajayj@berkeley.edu

Peter Abbeel
UC Berkeley
pabbeel@cs.berkeley.edu

Abstract

We present high quality image synthesis results using diffusion probabilistic models, a class of latent variable models inspired by considerations from nonequilibrium thermodynamics. Our models are trained by training on a weighted variational bound designed according to a novel connection between score-based probabilistic models and denoising score matching with Langevin dynamics, and our models naturally admit a progressive lossy decomposition scheme that can be interpreted as a generalization of progressive denoising. On the CIFAR-10 dataset, we obtain an inception score of 9.46 and a state-of-the-art FID of 2.17. On 256x256 LSUN, we obtain sample quality similar to ProgressiveGAN. Our implementation is available at <https://github.com/jonathanho/diffusion>.

1 Introduction

Generative models of all kinds have recently exhibited high quality samples in a wide variety of data modalities. Generative adversarial networks (GANs), autoregressive models, flows, and variational auto-encoders (VAEs) have all shown strong image generation results [27, 5, 58, 38, 25, 10, 32, 44, 57, 26, 33, 45], and there have been remarkable advances in energy-based modeling and score matching that have produced images comparable to those of GANs [11, 55].



Figure 1: Generated samples on CelebA-HQ 256 x 256 (left) and unconditional CIFAR10 (right)

34th Conference on Neural Information Processing Systems (NeurIPS 2020), Vancouver, Canada.

arXiv:2006.11239v2 [cs.LG] 16 Dec 2020

Published as a conference paper at ICLR 2021

SCORE-BASED GENERATIVE MODELING THROUGH STOCHASTIC DIFFERENTIAL EQUATIONS

Yang Song^{*}
Stanford University
yangsong@cs.stanford.edu

Jascha Sohl-Dickstein
Google Brain
jaschaid@google.com

Diederik P. Kingma
Google Brain
durk@google.com

Abhishek Kumar
Google Brain
abhishek@google.com

Stefano Ermon
Stanford University
ermon@cs.stanford.edu

Ben Poole
Google Brain
poole@google.com

ABSTRACT

Creating noise from data is easy; creating data from noise is generative modeling. We present a stochastic differential equation (SDE) that smoothly transforms a complex data distribution to a known prior distribution by slowly injecting noise, and a corresponding generative model that can sample from the target distribution by reversing the data distribution by slowly removing the noise. Crucially, the reversed-time SDE depends only on the time-dependent gradient field (a.k.a., score) of the perturbed data distribution. By leveraging advances in score-based generative modeling, we can accurately estimate models with millions of parameters and solve the resulting SDE solvers to generate samples. We show that this framework encapsulates previous approaches in score-based generative modeling: diffusion probabilistic modeling, allowing for new sampling procedures and new modeling capabilities. In particular, we show a principled connection between the two extremes in the evolution of the discretized reverse-time SDE. We also derive an equivalent neural ODE that samples from the same distribution as the SDE, but additionally enables exact likelihood computation, and improved sampling efficiency. In addition, we provide a new way to solve inverse problems with generative models, as demonstrated on examples of image denoising, generation, image inpainting, and colorization. Combined with multiple architectural improvements, we achieve record-breaking performance for unconditional image generation on CIFAR-10 with an inception score of 9.89 and FID of 2.20, a competitive likelihood of 2.99 bits/dim, and demonstrate high fidelity generation of 1024×1024 images for the first time from a score-based generative model.

1 INTRODUCTION

Two successful classes of probabilistic generative models involve sequentially corrupting training data with slowly increasing noise, and then learning to reverse the corruption in order to form a generative model of the data. *Score matching with Langevin dynamics* (SMLD) (Song & Ermon, 2019) estimates the gradient of the log probability with respect to data at each noise scale, and then uses Langevin dynamics to sample from a sequence of decreasing noise scales during generation. *Denoising diffusion probabilistic modeling* (DDPM) (Sohl-Dickstein et al., 2015; Ho et al., 2020) trains a sequence of probabilistic models to reverse each step of the noise corruption, using knowledge of the evolution of the reverse distribution to make training tractable. For example, we propose the DDPM training objective to directly compute scores at coarse noise scale. We therefore refer to these two model classes together as *score-based generative models*.

Score-based generative models, and related techniques (Bordes et al., 2017; Goyal et al., 2019; Du & Mordatch, 2019), have proven effective at generation of images (Song & Ermon, 2019; 2020; Ho et al., 2020), audio (Chen et al., 2020; Kong et al., 2020), graphics (Niu et al., 2020), and shapes (Cai et al., 2020).

^{*}Work partially done during an internship at Google Brain.

1

arXiv:2011.13456v2 [cs.LG] 10 Feb 2021

Score-based/Diffusion Generative Models

Score-based generative models, and diffusion generative models:

Deep Unsupervised Learning using Nonequilibrium Thermodynamics
Stanford University
JASCHA@STANFORD.EDU

Eric A. Weiss
University of California, Berkeley
EWEISS@BERKELEY.EDU

Niru Maheswaranathan
Stanford University
NIRUM@STANFORD.EDU

Surya Ganguli
Stanford University
SGANGULI@STANFORD.EDU

Abstract
A central problem in machine learning involves modeling complex data-sets using highly flexible models of the underlying distribution. Data, learning, sampling, inference, and evaluation are still analytically or computationally tractable. Here, we develop an approach that simultaneously achieves both flexibility and tractability. The essential idea, inspired by non-equilibrium statistical mechanics, is to learn and then destroy structure in a data distribution through an iterative forward diffusion process. We then learn a reverse diffusion process that restores structure in data, yielding a highly flexible and tractable generative model of the data. This approach allows us to learn, sample, generate, and evaluate probabilities in deep generative models with thousands of layers or time steps, as well as to compute conditional and posterior probabilities under the learned model. We additionally release an open source reference implementation of the algorithm.

1. Introduction
Historically, probabilistic models suffer from a tradeoff between two conflicting objectives: *tractability* and *flexibility*. Models that are *tractable* can be analytically evaluated and easily fit to data (e.g. a Gaussian or Laplace). However, non-parametric methods can be seen as transitioning smoothly between tractable and flexible models. For instance, a non-parametric Gaussian mixture model will represent a small amount of data using a single Gaussian, but may represent infinite data as a mixture of an infinite number of Gaussians.

arXiv:1503.03585v8 [cs.LG] 18 Nov 2015

Generative Modeling by Estimating Gradients of the Data Distribution
Stanford University
Yang Song
yongsong@cs.stanford.edu
UC Berkeley
Stefano Ermon
ermon@cs.stanford.edu

Abstract
We introduce a new generative model where samples are produced via Langevin dynamics using gradients of the data distribution estimated with score matching. Because gradients can be ill-defined and hard to estimate when the data resides on low-dimensional manifolds, we perform the score estimation with two Gaussian noise, and jointly estimate the corresponding score (i.e., the vector fields of gradients of the perturbed data distribution for all noise levels). For sampling, we propose an annealed Langevin dynamics where we use gradients corresponding to gradually decreasing noise levels in the score estimation process to close the loop. Our framework allows flexible model architectures, requires no sampling during training or the use of adversarial methods, and provides a learning objective that can be used for principled model comparisons. Our models produce samples comparable to GANs on MNIST, CelebA, and CIFAR-10 datasets, achieving a new state-of-the-art inception score of 8.87 on CIFAR-10. Additionally, we demonstrate that our models learn effective representations via image inpainting experiments.

1. Introduction
Generative models have many applications in machine learning. To list a few, they have been used to generate high-fidelity images [26, 6], synthesize realistic speech and music fragments [58], improve the performance of semi-supervised learning [27], generate synthetic samples and other anomalous data [54], improve learning [22], and explore promising states in reinforcement learning [41]. Recent progress is mainly driven by two approaches: likelihood-based methods [17, 29, 11, 60] and generative adversarial networks (GAN) [15]. The former uses log-likelihood (or a suitable surrogate) as the training objective, while the latter uses adversarial training to minimize $J(\theta)$ and $J(\phi)$ [15].

Although likelihood-based models and GANs have achieved great success, they have some intrinsic limitations. For example, likelihood-based models either have to use specialized architectures to build a normalized probability model (e.g., autoregressive models, flow models), or use surrogate losses (e.g., the evidence lower bound used in variational auto-encoders [29], contrastive divergence in generative models [11]), for training. GANs are some of the illustrations of likelihood-based models, but they are not directly comparable to likelihood-based generative models. In addition, the GAN objective is not suitable for evaluating and comparing different GAN models. While other objectives exist for generative modeling, such as noise contrastive estimation [19] and minimum probability flow [50], these methods typically only work well for low-dimensional data.

In this paper, we explore a new paradigm for generative modeling based on estimating gradients from the log-density [33] of the logarithmic data density, which is the gradient of the log-density function at the input data point. This is a vector field pointing in the direction where the log data density grows the most. We use a neural network trained with score matching [24] to learn this vector field from data. We then produce samples using Langevin dynamics, which approximately

arXiv:1907.05600v3 [cs.LG] 10 Oct 2020

Denoising Diffusion Probabilistic Models
Stanford University
Jonathan Ho
jonathanho@berkeley.edu
UC Berkeley
Ajay Jain
ajayj@berkeley.edu
UC Berkeley
Pieter Abbeel
pabbeel@cs.berkeley.edu

Abstract
We present high quality image synthesis results using diffusion probabilistic models, a class of latent variable models inspired by considerations from nonequilibrium thermodynamics. Our models are trained by training on a weighted variational bound designed according to a novel connection between score-based probabilistic models and denoising score matching with Langevin dynamics, and our models naturally admit a progressive lossy decomposition scheme that can be interpreted as a generalization of progressive denoising. On the CIFAR-10 dataset, we obtain an inception score of 9.46 and a state-of-the-art FID of 2.17. On 256x256 LSUN, we obtain sample quality similar to ProgressiveGAN. Our implementation is available at <https://github.com/jonathanho/diffusion>.

1. Introduction
Generative models of all kinds have recently exhibited high quality samples in a wide variety of data modalities. Generative adversarial networks (GANs), autoregressive models, flows, and variational auto-encoders (VAEs) have all shown strong image generation results [27, 5, 58, 38, 25, 10, 32, 44, 57, 26, 33, 45], and there have been remarkable advances in energy-based modeling and score matching that have produced images comparable to those of GANs [11, 55].



Figure 1: Generated samples on CelebA-HQ 256 x 256 (left) and unconditional CIFAR10 (right)

arXiv:2006.11239v2 [cs.LG] 16 Dec 2020

SCORE-BASED GENERATIVE MODELING THROUGH STOCHASTIC DIFFERENTIAL EQUATIONS
Stanford University
Yang Song*
yangsong@cs.stanford.edu
Google Brain
Jascha Sohl-Dickstein
jaschad@google.com
Google Brain
Diederik P. Kingma
dirk@google.com

Abhishek Kumar
Google Brain
abhishek@google.com
Stanford University
Stefano Ermon
ermon@cs.stanford.edu
Google Brain
Ben Poole
Google Brain
poole@google.com

ABSTRACT
Creating noise from data is easy; creating data from noise is generative modeling. We present a stochastic differential equation (SDE) that smoothly transforms a complex data distribution to a known prior distribution by slowly injecting noise, and a corresponding generative model that can be trained to sample from the target data distribution by slowly removing the noise. Crucially, the reverse-time SDE depends only on the time-dependent gradient field (a.k.a., score) of the perturbed data distribution. By leveraging advances in score-based generative modeling, we can accurately estimate samples with a sequence of denoising steps using standard SDE solvers to generate samples. We show that this framework encapsulates previous approaches in score-based generative modeling: diffusion probabilistic modeling, allowing for new sampling procedures and new modeling capabilities. In particular, we show a principled connection between the noise corruption in the evolution of the discretized reverse-time SDE. We also derive an equivalent neural ODE that samples from the same distribution as the SDE, but additionally enables exact likelihood computation, and improved sampling efficiency. In addition, we provide a new way to solve inverse problems with generative models, as demonstrated on examples of image classification, generation, image inpainting, and colorization. Combined with multiple architectural improvements, we achieve record-breaking performance for unconditional image generation on CIFAR-10 with an inception score of 9.89 and FID of 2.20, a competitive likelihood of 2.99 bits/dim, and demonstrate high fidelity generation of 1024 x 1024 images for the first time from a score-based generative model.

1. INTRODUCTION
Two successful classes of probabilistic generative models involve sequentially corrupting training data with slowly increasing noise, and then learning to reverse the corruption in order to form a generative model of the data. *Score matching with Langevin dynamics* (SMLD) (Song & Ermon, 2019) estimates the probability of log-probability with respect to data at each noise scale, and then uses Langevin dynamics to sample from a sequence of decreasing noise scales during generation. *Denoising diffusion probabilistic modeling* (DDPM) (Sohl-Dickstein et al., 2015; Ho et al., 2020) trains a sequence of probabilistic models to reverse each step of the noise corruption, using knowledge of the evolution of the reverse distribution to make training tractable. For this paper, we expand the DDPM framework to generatively corrupt scores at coarse scales. We therefore refer to these two model classes together as *score-based generative models*.

Score-based generative models, and related techniques (Bordes et al., 2017; Goyal et al., 2019; Du & Mordatch, 2019), have proven effective at generation of images (Song & Ermon, 2019; 2020; Ho et al., 2020), audio (Chen et al., 2020; Kong et al., 2020), shapes (Giu et al., 2020), and

*Work partially done during an internship at Google Brain.

1

33rd Conference on Neural Information Processing Systems (NeurIPS 2019), Vancouver, Canada.
34th Conference on Neural Information Processing Systems (NeurIPS 2020), Vancouver, Canada.

First Diffusion Model (ICML 2015)

Score-based/Diffusion Generative Models

Score-based generative models, and diffusion generative models:

arXiv:1503.03585v8 [cs.LG] 18 Nov 2015

Deep Unsupervised Learning using Nonequilibrium Thermodynamics

Jascha Sohl-Dickstein
Stanford University

Eric A. Weiss
University of California, Berkeley

Niru Maheswaranathan
Stanford University

Surya Ganguli
Stanford University

Abstract

A central problem in machine learning involves modeling complex data-sets using highly flexible models of the underlying distribution. Data fitting, sampling, inference, and evaluation are still analytically or computationally tractable. Here, we develop an approach that simultaneously achieves both flexibility and tractability. The essential idea, inspired by non-equilibrium statistical mechanics, is to learn and then destroy structure in a data distribution through an iterative forward diffusion process. We then learn a reverse diffusion process that restores structure in data, yielding a highly flexible and tractable generative model of the data. This approach allows us to learn, sample, infer, and evaluate probabilities in deep generative models with thousands of layers or time steps, as well as to compute conditional and posterior probabilities under the learned model. We additionally release an open source reference implementation of the algorithm.

1. Introduction

Historically, probabilistic models suffer from a tradeoff between two conflicting objectives: *tractability* and *flexibility*. Models that are *tractable* can be analytically evaluated and easily fit to data (e.g. a Gaussian or Laplace). However,

variational Bayes (Jordan et al., 1999), contrastive divergence (Welling & Hinton, 2002; Hinton, 2002), minimum probability flow (Sohl-Dickstein et al., 2011a), variational auto-encoder (Kingma & Welling, 2011), score matching (Gretton & Sugiyama, 2009; Gretton et al., 2012), maximum likelihood (Hyvonen, 2005), pseudolikelihood (Besag, 1975), keep belief propagation (Murphy et al., 1999), and many, many more, Non-parametric methods (Gershman & Blei, 2012) can also be very effective¹.

- 1.1. Diffusion probabilistic models
- We present a novel way to define probabilistic models that allows:

1. extreme flexibility in model structure,
2. exact sampling.

¹Non-parametric methods can be seen as transitioning smoothly between tractable and flexible models. For instance, a non-parametric Gaussian mixture model will represent a small amount of data using a single Gaussian, but may represent infinite data as a mixture of an infinite number of Gaussians.

Score-based Generative Model (NeurIPS 2019)

arXiv:1907.05600v3 [cs.LG] 10 Oct 2020

Generative Modeling by Estimating Gradients of the Data Distribution

Yang Song
Stanford University
yangsong@cs.stanford.edu

Stefano Ermon
Stanford University
ermon@cs.stanford.edu

Abstract

We introduce a new generative model where samples are produced via Langevin dynamics using gradients of the data distribution estimated with score matching. Because gradients can be ill-defined and hard to estimate when the data resides on low-dimensional manifolds, we perform the data-to-distribution warping of Gaußians, noise, and jointly estimate the corresponding score (i.e., the vector fields of gradients of the perturbed data distribution for all noise levels). For sampling, we propose an annealed Langevin dynamics where we use gradients corresponding to gradually increasing noise levels in the annealing process to close the loop. The final One. Our framework allows flexible model architectures, requires no sampling during training or the use of adversarial methods, and provides a learning objective that can be used for principled model comparisons. Our models produce samples comparable to GANs on MNIST, CelebA, and CIFAR-10 datasets, achieving a new state-of-the-art inception score of 8.87 on CIFAR-10. Additionally, we demonstrate that our models learn effective representations via image inpainting experiments.

1 Introduction

Generative models have many applications in machine learning. To list a few, they have been used to generate high-fidelity images [26, 6], synthesize realistic speech and music fragments [58], improve the performance of semi-supervised learning [20], generate synthetic samples and other anomalous data [54], improve learning [22], and explore promising states in reinforcement learning [41]. Recent progress is mainly driven by two approaches: likelihood-based methods [17, 29, 11, 60] and generative adversarial networks (GAN [15]). The former uses log-likelihood (or a suitable surrogate) as the training objective, while the latter uses adversarial training to minimize $J(\theta, \phi)$ [15, 20, 21].

Although likelihood-based models and GANs have achieved great success, they have some intrinsic limitations.

For example, likelihood-based models either have to use specialized architectures to build a normalized probability model (e.g., autoregressive models, flow models), or use surrogate losses (e.g., the evidence lower bound used in variational auto-encoders [29]), contrastive divergence in generative models [21] for training. GANs are some of the manifestations of likelihood-based models, but they are not necessarily able to learn a sensible generative process. In addition, the GAN objective is not suitable for evaluating and comparing different GAN models. While other objectives exist for generative modeling, such as noise contrastive estimation [19] and minimum probability flow [50], these methods typically only work well for low-dimensional data.

In this paper, we explore a new paradigm for generative modeling based on estimating gradients from the log-density [33] of the logarithmic data density, which is the gradient of the log-density function at the input data point. This is a vector field pointing in the direction where the log data density grows the most. We use a neural network trained with score matching [24] to learn this vector field from data. We then produce samples using Langevin dynamics, which approximately

33rd Conference on Neural Information Processing Systems (NeurIPS 2019), Vancouver, Canada.

arXiv:2006.11239v2 [cs.LG] 16 Dec 2020

Denoising Diffusion Probabilistic Models

Jonathan Ho
UC Berkeley
jonathanho@berkeley.edu

Ajay Jain
UC Berkeley
ajayj@berkeley.edu

Pieter Abbeel
UC Berkeley
pabbeel@cs.berkeley.edu

Abstract

We present high quality image synthesis results using diffusion probabilistic models, a class of latent variable models inspired by considerations from nonequilibrium thermodynamics. Our models are trained by training on a weighted empirical loss function designed according to a causal connection between denoising probabilistic models and denoising score matching with Langevin dynamics, and our models naturally admit a progressive losy decomposition scheme that can be interpreted as a generalization of progressive denoising. On the CIFAR-10 dataset, we obtain an inception score of 9.46 and a state-of-the-art FID of 2.17. On 256x256 LSUN, we obtain sample quality similar to ProgressiveGAN. Our implementation is available at <https://github.com/jonathanho/diffusion>.

1 Introduction

Generative models of all kinds have recently exhibited high quality samples in a wide variety of data modalities. Generative adversarial networks (GANs), autoregressive models, flows, and variational auto-encoders (VAEs) have all shown strong image generation results [27, 5, 58, 38, 25, 10, 32, 44, 57, 26, 33, 45], and there have been remarkable advances in energy-based modeling and score matching that have produced images comparable to those of GANs [11, 55].



Figure 1: Generated samples on CelebA-HQ 256 × 256 (left) and unconditional CIFAR10 (right)

34th Conference on Neural Information Processing Systems (NeurIPS 2020), Vancouver, Canada.

Published as a conference paper at ICLR 2021

SCORE-BASED GENERATIVE MODELING THROUGH STOCHASTIC DIFFERENTIAL EQUATIONS

Yang Song^{*}
Stanford University
yangsong@cs.stanford.edu

Jascha Sohl-Dickstein
Google Brain
jaschaa@google.com

Abhishek Kumar
Google Brain
abhishek@google.com

Stefano Ermon
Stanford University
ermon@cs.stanford.edu

Ben Poole
Google Brain
poole@google.com

ABSTRACT

Creating noise from data is easy; creating data from noise is generative modeling. We present a stochastic differential equation (SDE) that smoothly transforms a complex data distribution to a known prior distribution by slowly injecting noise, and a corresponding reverse SDE that smoothly transforms a known prior distribution back into data distribution by slowly removing the noise. Crucially, the reverse-time SDE depends only on the time-dependent gradient field (a.k.a., score) of the perturbed data distribution. By leveraging advances in score-based generative modeling, we can accurately estimate models with millions of parameters to generate samples using SDE solvers to generate samples. We show that this framework encapsulates previous approaches in score-based generative modeling: diffusion probabilistic modeling, allowing for new sampling procedures and new modeling capabilities. In particular, we show a principled connection between the noise corruption in the evolution of the discretized reverse-time SDE. We also derive an equivalent neural ODE that samples from the same distribution as the SDE, but additionally enables exact likelihood computation, and improved sampling efficiency. In addition, we provide a new way to solve inverse problems with generative models, as demonstrated on examples of image classification, generation, image inpainting, and colorization. Combined with multiple architectural improvements, we achieve record-breaking performance for unconditional image generation on CIFAR-10 with an inception score of 9.89 and FID of 2.20, a competitive likelihood of 2.99 bits/dim, and demonstrate high fidelity generation of 1024×1024 images for the first time from a score-based generative model.

1 INTRODUCTION

Two successful classes of probabilistic generative models involve sequentially corrupting training data with slowly increasing noise, and then learning to reverse the corruption in order to form a generative model of the data. *Score matching with Langevin dynamics* (SMD) (Song & Ermon, 2019) approximates the log-density of the log-probability with respect to the data at each noise scale, and then uses Langevin dynamics to sample from a sequence of noise corrupted scales during generation. *Denoising diffusion probabilistic modeling* (DDPM) (Sohl-Dickstein et al., 2015; Ho et al., 2020) trains a sequence of probabilistic models to reverse each step of the noise corruption, using knowledge of the evolution of the reverse distribution to make training tractable. For example, we expand the DDPM training objective to include scores at coarse noise scale. We therefore refer to these two model classes together as *score-based generative models*.

Score-based generative models, and related techniques (Bordes et al., 2017; Goyal et al., 2017; Du & Mordatch, 2019), have proven effective at generation of images (Song & Ermon, 2019; 2020; Ho et al., 2020), audio (Chen et al., 2020; Kong et al., 2020), graphics (Niu et al., 2020), and shapes (Cai et al., 2020).

^{*}Work partially done during an internship at Google Brain.

Score-based/Diffusion Generative Models

Score-based generative models, and diffusion generative models:

arXiv:1503.03585v8 [cs.LG] 18 Nov 2015

Deep Unsupervised Learning using Nonequilibrium Thermodynamics

Jascha Sohl-Dickstein
Stanford University

Eric A. Weiss
University of California, Berkeley

Niru Maheswaranathan
Stanford University

Surya Ganguli
Stanford University

Abstract

A central problem in machine learning involves modeling complex data-sets using highly flexible models of the underlying distribution. Data fitting, sampling, inference, and evaluation are still analytically or computationally tractable. Here, we develop an approach that simultaneously achieves both flexibility and tractability. The essential idea, inspired by non-equilibrium statistical mechanics, is to learn and then destroy structure in a data distribution through an iterative forward diffusion process. We then learn a reverse diffusion process that restores structure in data, yielding a highly flexible and tractable generative model of the data. This approach allows us to learn, sample, infer, and evaluate probabilities in deep generative models with thousands of layers or time steps, as well as to compute conditional and posterior probabilities under the learned model. We additionally release an open source reference implementation of the algorithm.

1. Introduction

Historically, probabilistic models suffer from a tradeoff between two conflicting objectives: *tractability* and *flexibility*. Models that are *tractable* can be analytically evaluated and easily fit to data (e.g. a Gaussian or Laplace). However, a central problem in machine learning involves modeling complex data-sets using highly flexible models of the underlying distribution. Data fitting, sampling, inference, and evaluation are still analytically or computationally tractable. Here, we develop an approach that simultaneously achieves both flexibility and tractability. The essential idea, inspired by non-equilibrium statistical mechanics, is to learn and then destroy structure in a data distribution through an iterative forward diffusion process. We then learn a reverse diffusion process that restores structure in data, yielding a highly flexible and tractable generative model of the data. This approach allows us to learn, sample, infer, and evaluate probabilities in deep generative models with thousands of layers or time steps, as well as to compute conditional and posterior probabilities under the learned model. We additionally release an open source reference implementation of the algorithm.

1.1. Diffusion probabilistic models

We present a novel way to define probabilistic models that allows:

1. extreme flexibility in model structure,
2. exact sampling.

¹Non-parametric methods can be seen as transitioning smoothly between tractable and flexible models. For instance, a non-parametric Gaussian mixture model will represent a small amount of data using a single Gaussian, but may represent infinite data as a mixture of an infinite number of Gaussians.

JASCHA @ STANFORD.EDU

EAWIESS @ BERKELEY.EDU

NIRUM @ STANFORD.EDU

SGANGULI @ STANFORD.EDU

arXiv:1907.05600v3 [cs.LG] 10 Oct 2020

arXiv:2006.11239v2 [cs.LG] 16 Dec 2020

Generative Modeling by Estimating Gradients of the Data Distribution

Yang Song
Stanford University
yangsong@cs.stanford.edu

Stefano Ermon
Stanford University
ermon@cs.stanford.edu

Abstract

We introduce a new generative model where samples are produced via Langevin dynamics using gradients of the data distribution estimated with score matching. Because gradients can be ill-defined and hard to estimate when data resides on low-dimensional manifolds, we perform the data-to-distribution warping (i.e., noise, and joint) using the corresponding score (i.e., the vector fields of gradients of the perturbed data distribution for all noise levels). For sampling, we propose an annealed Langevin dynamics where we use gradients corresponding to gradually increasing noise levels in the annealing process to close the loop. One benefit. Our framework allows flexible model architectures, reduces no sampling during training or the use of adversarial methods, and provides a learning objective that can be used for principled model comparisons. Our models produce samples comparable to GANs on MNIST, CelebA, and CIFAR-10 datasets, achieving a new state-of-the-art inception score of 8.87 on CIFAR-10. Additionally, we demonstrate that our models learn effective representations via image inpainting experiments.

1 Introduction

Generative models have many applications in machine learning. To list a few, they have been used to generate high-fidelity images [26, 6], synthesize realistic speech and music fragments [58], improve the performance of semi-supervised learning [20], generate realistic samples and other anomalous data [54], image learning [22], and explore promising states in reinforcement learning [41]. Recent progress is mainly driven by two approaches: likelihood-based methods [17, 29, 11, 60] and generative adversarial networks (GAN [15]). The former uses log-likelihood (or a suitable surrogate) as the training objective, while the latter uses adversarial training to minimize $J(\theta)$ and $J(\phi)$ [15].

Although likelihood-based models and GANs have achieved great success, they have some intrinsic limitations. For example, likelihood-based models either have to use specialized architectures to build a normalized probability model (e.g., autoregressive models, flow models), or use a surrogate loss (e.g., the evidence lower bound used in variational auto-encoders [29]), contrastive divergence in energy-based models [21] for training. GANs are some of the illustrations of likelihood-based models, but they are not suitable for evaluating and comparing different GAN models. In addition, the GAN objective is not suitable for evaluating and comparing different GAN models. While other objectives exist for generative modeling, such as noise contrastive estimation [19] and minimum probability flow [50], these methods typically only work well for low-dimensional data.

In this paper, we explore a new paradigm for generative modeling based on estimating gradients from the log-density [33] of the logarithmic data density, which is the gradient of the log-density function at the input data point. This is a vector field pointing in the direction where the log data density grows the most. We use a neural network trained with score matching [24] to learn this vector field from data. We then produce samples using Langevin dynamics, which approximately

jonathanho@berkeley.edu

Ajay Jain
UC Berkeley
ajayj@berkeley.edu

Pieter Abbeel
UC Berkeley
pabbeel@cs.berkeley.edu

Abstract

We present high quality image synthesis using diffusion probabilistic models, a class of latent variable models inspired by considerations from nonequilibrium thermodynamics. Our models are trained by training on a weighted variational bound designed according to a causal connection between score-based probabilistic models and denoising score matching with Langevin dynamics, and our models naturally admit a progressive lossy decomposition scheme that can be interpreted as a generalization of progressive denoising and a state-of-the-art FID of 9.47. On 256x256 LSUN, we obtain sample quality similar to ProgressiveGAN. Our implementation is available at <https://github.com/jonathanho/diffusion>.

arXiv:2011.13456v2 [cs.LG] 10 Feb 2021

arXiv:2006.11239v2 [cs.LG] 16 Dec 2020

Denoising Diffusion Probabilistic Models

Jonathan Ho
UC Berkeley
jonathanho@berkeley.edu

Ajay Jain
UC Berkeley
ajayj@berkeley.edu

Pieter Abbeel
UC Berkeley
pabbeel@cs.berkeley.edu

Abstract

We present high quality image synthesis using diffusion probabilistic models, a class of latent variable models inspired by considerations from nonequilibrium thermodynamics. Our models are trained by training on a weighted variational bound designed according to a causal connection between score-based probabilistic models and denoising score matching with Langevin dynamics, and our models naturally admit a progressive lossy decomposition scheme that can be interpreted as a generalization of progressive denoising and a state-of-the-art FID of 9.47. On 256x256 LSUN, we obtain sample quality similar to ProgressiveGAN. Our implementation is available at <https://github.com/jonathanho/diffusion>.

1 Introduction

Generative models of all kinds have recently exhibited high quality samples in a wide variety of data modalities. Generative adversarial networks (GANs), autoregressive models, flows, and variational auto-encoders (VAEs) have all shown strong image generation results [27, 3, 58, 38, 25, 10, 32, 44, 57, 26, 33, 45], and there have been remarkable advances in energy-based modeling and score matching that have produced images comparable to those of GANs [11, 55].



Figure 1: Generated samples on CelebA-HQ 256 × 256 (left) and unconditional CIFAR10 (right)

34th Conference on Neural Information Processing Systems (NeurIPS 2020), Vancouver, Canada.

Published as a conference paper at ICLR 2021

SCORE-BASED GENERATIVE MODELING THROUGH STOCHASTIC DIFFERENTIAL EQUATIONS

Yang Song^{*}
Stanford University
yangsong@cs.stanford.edu

Jascha Sohl-Dickstein
Google Brain
jaschaa@google.com

Diederik P. Kingma
Google Brain
durk@google.com

Abhishek Kumar
Google Brain
abhishek@google.com

Stefano Ermon
Stanford University
ermon@cs.stanford.edu

Ben Poole
Google Brain
poole@google.com

ABSTRACT

Creating noise from data is easy; creating data from noise is generative modeling. We present a stochastic differential equation (SDE) that smoothly transforms a complex data distribution to a known prior distribution by slowly injecting noise, and a corresponding generative model that samples from the prior. Our SDE starts with a complex data distribution by slowly removing the noise. Crucially, the reverse-time SDE depends only on the time-dependent gradient field (a.k.a., score) of the perturbed data distribution. By leveraging advances in score-based generative modeling, we can accurately estimate the score function with a neural network and numerically solve SDE solvers to generate samples. We show that this framework encapsulates previous approaches in score-based generative modeling: diffusion probabilistic modeling, allowing for new sampling procedures and new modeling capabilities. In particular, we show a principled connection between the noise corruption in the evolution of the discretized reverse-time SDE. We also derive an equivalent neural ODE that samples from the same distribution as the SDE, but additionally enables exact likelihood computation, and improved sampling efficiency. In addition, we provide a new way to solve inverse problems with generative models, as demonstrated on examples of image classification, generation, image inpainting, and colorization. Combined with multiple architectural improvements, we achieve record-breaking performance for unconditional image generation on CIFAR-10 with an inception score of 9.89 and FID of 2.20, a competitive likelihood of 2.99 bits/dim, and demonstrate high fidelity generation of 1024 × 1024 images for the first time from a score-based generative model.

1 INTRODUCTION

Two successful classes of probabilistic generative models involve sequentially corrupting training data with slowly increasing noise, and then learning to reverse the corruption in order to form a generative model of the data. *Score matching with Langevin dynamics* (SMLD) (Song & Ermon, 2019) performs the product of the log probability with respect to the data at each noise scale, and then uses Langevin dynamics to sample from a sequence of decreasing noise scales during generation. *Denoising diffusion probabilistic modeling* (DDPM) (Sohl-Dickstein et al., 2015; Ho et al., 2020) trains a sequence of probabilistic models to reverse each step of the noise corruption, using knowledge of the gradient form of the reverse distribution to make training tractable. For example, we expand the DDPM training objective to include score at coarse noise scale. We therefore refer to these two model classes together as *score-based generative models*.

Score-based generative models, and related techniques (Bordes et al., 2017; Goyal et al., 2017; Du & Mordatch, 2019), have proven effective at generation of images (Song & Ermon, 2019; 2020; Ho et al., 2020), audio (Chen et al., 2020; Kong et al., 2020), shapes (Niu et al., 2020), and

^{*}Work partially done during an internship at Google Brain.

Denoising Diffusion Model (NeurIPS 2020)

Score-based/Diffusion Generative Models

Score-based generative models, and diffusion generative models:

arXiv:1503.03585v8 [cs.LG] 18 Nov 2015

Deep Unsupervised Learning using Nonequilibrium Thermodynamics

Jascha Sohl-Dickstein
Stanford University

Eric A. Weiss
University of California, Berkeley

Niru Maheswaranathan
Stanford University

Surya Ganguli
Stanford University

Abstract

A central problem in machine learning involves modeling complex data-sets using highly flexible models of the underlying distribution. Data fitting, sampling, inference, and evaluation are still analytically or computationally tractable. Here, we develop an approach that simultaneously achieves both flexibility and tractability. The essential idea, inspired by non-equilibrium statistical mechanics, is to learn and then destroy structure in a data distribution through an iterative forward diffusion process. We then learn a reverse diffusion process that restores structure in data, yielding a highly flexible and tractable generative model of the data. This approach allows us to learn, sample, infer, and evaluate probabilities in deep generative models with thousands of layers or time steps, as well as to compute conditional and posterior probabilities under the learned model. We additionally release an open source reference implementation of the algorithm.

1. Introduction

Historically, probabilistic models suffer from a tradeoff between two conflicting objectives: *tractability* and *flexibility*. Models that are *tractable* can be analytically evaluated and easily fit to data (e.g. a Gaussian or Laplace). However, non-parametric methods (Gershman & Blei, 2012) can also be very effective¹.

1.1. Diffusion probabilistic models

We present a novel way to define probabilistic models that allows:

1. extreme flexibility in model structure,
2. exact sampling.

¹Non-parametric methods can be seen as transitioning smoothly between tractable and flexible models. For instance, a non-parametric Gaussian mixture model will represent a small amount of data using a single Gaussian, but may represent infinite data as a mixture of an infinite number of Gaussians.

JASCHA @ STANFORD.EDU

EAWIESS @ BERKELEY.EDU

NIRUM @ STANFORD.EDU

SGANGULI @ STANFORD.EDU

arXiv:1907.05600v3 [cs.LG] 10 Oct 2020

arXiv:2006.11239v2 [cs.LG] 16 Dec 2020

Generative Modeling by Estimating Gradients of the Data Distribution

Yang Song
Stanford University
yangsong@cs.stanford.edu

Stefano Ermon
Stanford University
ermon@cs.stanford.edu

Abstract

We introduce a new generative model where samples are produced via Langevin dynamics using gradients of the data distribution estimated with score matching. Because gradients can be ill-defined and hard to estimate when the data resides on low-dimensional manifolds, we perform the data-to-distribution warping of Givens noise, and jointly estimate the corresponding scores (i.e., the vector fields of gradients of the perturbed data distribution for all noise levels). For sampling, we propose an annealed Langevin dynamics where we use gradients corresponding to gradually increasing noise levels in the score matching process to close the loop. Our framework allows flexible model architectures, requires no sampling during training or the use of adversarial methods, and provides a learning objective that can be used for principled model comparisons. Our models produce samples comparable to GANs on MNIST, CelebA, and CIFAR-10 datasets, achieving a new state-of-the-art inception score of 8.87 on CIFAR-10. Additionally, we demonstrate that our models learn effective representations via image inpainting experiments.

1 Introduction

Generative models have many applications in machine learning. To list a few, they have been used to generate high-fidelity images [26, 6], synthesize realistic speech and music fragments [58], improve the performance of semi-supervised learning [27], generate realistic samples and other anomalous data [54], improve learning [22], and explore promising states in reinforcement learning [41]. Recent progress is mainly driven by two approaches: likelihood-based methods [17, 29, 11, 60] and generative adversarial networks (GAN [15]). The former uses log-likelihood (or a suitable surrogate) as the training objective, while the latter uses adversarial training to minimize $J(\theta)$ and $J(\phi)$ [15].

Although likelihood-based models and GANs have achieved great success, they have some intrinsic limitations. For example, likelihood-based models either have to use specialized architectures to build a normalized probability model (e.g., autoregressive models, flow models), or use surrogate losses (e.g., the evidence lower bound used in variational auto-encoders [29]), contrastive divergence in energy-based models [21] for training. GANs are some of the illustrations of likelihood-based models, but they are not suitable for evaluating and comparing different GAN models. While other objectives exist for generative modeling, such as noise contrastive estimation [19] and minimum probability flow [50], these methods typically only work well for low-dimensional data.

In this paper, we explore a new paradigm for generative modeling based on estimating gradients from the log-density [33] of the logarithmic data density, which is the gradient of the log data density at the input data point. This is a vector field pointing in the direction where the log data density grows the most. We use a neural network trained with score matching [24] to learn this vector field from data. We then produce samples using Langevin dynamics, which approximately

33rd Conference on Neural Information Processing Systems (NeurIPS 2019), Vancouver, Canada.

arXiv:2011.13456v2 [cs.LG] 10 Feb 2021

Denoising Diffusion Probabilistic Models

Jonathan Ho
UC Berkeley
jonathanho@berkeley.edu

Ajay Jain
UC Berkeley
ajayj@berkeley.edu

Peter Abbeel
UC Berkeley
pabbeel@cs.berkeley.edu

Abstract

We present high quality image synthesis results using diffusion probabilistic models, a class of latent variable models inspired by considerations from nonequilibrium thermodynamics. Our models are trained by training on a weighted empirical loss function designed according to a causal connection between score-based probabilistic models and denoising score matching with Langevin dynamics, and our models naturally admit a progressive losy decompression scheme that can be interpreted as a generalization of progressive denoising and a state-of-the-art FID on CIFAR10 dataset. We obtain an inception score of 9.46 and a state-of-the-art FID of 2.17. On 256x256 LSUN, we obtain sample quality similar to ProgressiveGAN. Our implementation is available at <https://github.com/jonathanho/diffusion>.

1 Introduction

Generative models of all kinds have recently exhibited high quality samples in a wide variety of data modalities. Generative adversarial networks (GANs), autoregressive models, flows, and variational auto-encoders (VAEs) have all shown strong image generation results [27, 3, 58, 38, 25, 10, 32, 44, 57, 26, 33, 45], and there have been remarkable advances in energy-based modeling and score matching that have produced images comparable to those of GANs [11, 55].



Figure 1: Generated samples on CelebA-HQ 256 x 256 (left) and unconditional CIFAR10 (right)

34th Conference on Neural Information Processing Systems (NeurIPS 2020), Vancouver, Canada.

Published as a conference paper at ICLR 2021

SCORE-BASED GENERATIVE MODELING THROUGH STOCHASTIC DIFFERENTIAL EQUATIONS

Yang Song^{*}
Stanford University
yangsong@cs.stanford.edu

Jascha Sohl-Dickstein
Google Brain
jaschaa@google.com

Diederik P. Kingma
Google Brain
durk@google.com

Abhishek Kumar
Google Brain
abhishek@google.com

Stefano Ermon
Stanford University
ermon@cs.stanford.edu

Ben Poole
Google Brain
poole@google.com

ABSTRACT

Creating noise from data is easy; creating data from noise is generative modeling. We present a stochastic differential equation (SDE) that smoothly transforms a complex data distribution to a known prior distribution by slowly injecting noise, and a corresponding generative model that can sample from the target distribution by reversing the data distribution by slowly removing the noise. Crucially, the reversed-time SDE depends only on the time-dependent gradient field (a.k.a., score) of the perturbed data distribution. By leveraging advances in score-based generative modeling, we can accurately estimate models with millions of parameters and efficiently use SDE solvers to generate samples. We show that this framework encapsulates previous approaches in score-based generative modeling: diffusion probabilistic modeling, allowing for new sampling procedures and new modeling capabilities. In particular, we show a principled connection between the noise corruption in the evolution of the discretized reverse-time SDE. We also derive an equivalent neural ODE that samples from the same distribution as the SDE, but additionally enables exact likelihood computation, and improved sampling efficiency. In addition, we provide a new way to solve inverse problems with generative models, as demonstrated on examples of image classification, generation, image inpainting, and colorization. Combined with multiple architectural improvements, we achieve record-breaking performance for unconditional image generation on CIFAR-10 with an inception score of 9.89 and FID of 2.20, a competitive likelihood of 2.99 bits/dim, and demonstrate high fidelity generation of 1024×1024 images for the first time from a score-based generative model.

1 INTRODUCTION

Two successful classes of probabilistic generative models involve sequentially corrupting training data with slowly increasing noise, and then learning to reverse the corruption in order to form a generative model of the data. *Score matching with Langevin dynamics* (SMLD) (Song & Ermon, 2019) performs the perturbation of semi-supervised data [29] and generates samples on a noise scale, and then uses Langevin dynamics to sample from a sequence of decreasing noise scales during generation. *Denoising diffusion probabilistic modeling* (DDPM) (Sohl-Dickstein et al., 2015; Ho et al., 2020) trains a sequence of probabilistic models to reverse each step of the noise corruption, using knowledge of the evolution of the reverse diffusion to make training tractable. For example, we expand the DDPM training objective to directly corrupt scores at coarse scales. We therefore refer to these two model classes together as *score-based generative models*.

Score-based generative models, and related techniques (Bordes et al., 2017; Goyal et al., 2017; Du & Mordatch, 2019), have proven effective at generation of images (Song & Ermon, 2019; 2020; Ho et al., 2020), audio (Chen et al., 2020; Kong et al., 2020), shapes (Cai et al., 2020), and text (Ho et al., 2020).

^{*}Work partially done during an internship at Google Brain.

Score-based Models via SDEs (ICLR 2021)

Score-based/Diffusion Generative Models

Score-based and diffusion generative models take advantage of the approximate inference algorithms we have recently learned!

Score-based/Diffusion Generative Models

Score-based and diffusion generative models take advantage of the approximate inference algorithms we have recently learned!

Score-based generative models → [today]

- Involve Langevin Monte Carlo (LMC) with an annealed proposal distribution.

Score-based/Diffusion Generative Models

Score-based and diffusion generative models take advantage of the approximate inference algorithms we have recently learned!

Score-based generative models → [today]

- Involve Langevin Monte Carlo (LMC) with an annealed proposal distribution.

Denoising diffusion generative models → [upcoming]

- Involve variational inference (VI) via ELBO maximization.

Score-based Generative Models – Setup

Score-based Generative Models – Setup

Suppose that we have a dataset: $D = \{x_1, x_2, \dots, x_n\}$

Where each data point is assumed to be drawn from some unknown distribution, i.e.,

$$x_i \sim p(x)$$

Score-based Generative Models – Setup

Suppose that we have a dataset: $D = \{x_1, x_2, \dots, x_n\}$

Where each data point is assumed to be drawn from some unknown distribution, i.e.,

$$x_i \sim p(x)$$

Suppose we want to model this distribution using a model $p_\theta(x)$, where we assume this distribution is parameterized by a $\theta \in \mathbb{R}^d$.

Score-based Generative Models – Setup

We can write this model $p_\theta(x)$ in a general way as:

Score-based Generative Models – Setup

We can write this model $p_\theta(x)$ in a general way as:

$$p_\theta(x) = \frac{1}{Z(\theta)} \tilde{p}_\theta(x) = \frac{1}{Z(\theta)} e^{-f_\theta(x)}$$

Where, as before,

Score-based Generative Models – Setup

We can write this model $p_\theta(x)$ in a general way as:

$$p_\theta(x) = \frac{1}{Z(\theta)} \tilde{p}_\theta(x) = \frac{1}{Z(\theta)} e^{-f_\theta(x)}$$

Where, as before, $Z(\theta)$ is the normalization constant.

Score-based Generative Models – Setup

We can write this model $p_\theta(x)$ in a general way as:

$$p_\theta(x) = \frac{1}{Z(\theta)} \tilde{p}_\theta(x) = \frac{1}{Z(\theta)} e^{-f_\theta(x)}$$

Where, as before, $Z(\theta)$ is the normalization constant.

And $\tilde{p}_\theta(x)$ is an unnormalized PDF.

Score-based Generative Models – Setup

We can write this model $p_\theta(x)$ in a general way as:

$$p_\theta(x) = \frac{1}{Z(\theta)} \tilde{p}_\theta(x) = \frac{1}{Z(\theta)} e^{-f_\theta(x)}$$

Where, as before, $Z(\theta)$ is the normalization constant.

And $\tilde{p}_\theta(x)$ is an unnormalized PDF.

And $f_\theta(x)$ is equal to the negative log of the unnormalized PDF, i.e.,

$$f_\theta(x) = -\log \tilde{p}_\theta(x)$$

Sometimes called an
energy-based model.

Score-based Generative Models – Difficulties with MLE

Previous strategies involve trying to learn parameters θ via maximum likelihood:

$$\theta^* = \operatorname{argmax}_{\theta} \sum_{i=1}^N \log p_{\theta}(x_i)$$

Score-based Generative Models – Difficulties with MLE

Previous strategies involve trying to learn parameters θ via maximum likelihood:

$$\theta^* = \operatorname{argmax}_{\theta} \sum_{i=1}^N \log p_{\theta}(x_i)$$

However:

- To optimize (or evaluate) $p_{\theta}(x)$, we need the normalization constant $Z(\theta)$.
- ...which is intractable in general for most $\tilde{p}_{\theta}(x)$.

Score-based Generative Models – Main Strategy

So instead, we will take the following strategy:

Score-based Generative Models – Main Strategy

So instead, we will take the following strategy:

- (1) **First**, we will aim to learn what is referred to as the **(Stein) score function** of the log data density.
- (2) **Second**, Given this score function, we will be able to draw samples from $p_\theta(x)$ using Langevin Monte Carlo.

Score-based Generative Models – Main Strategy

So instead, we will take the following strategy:

- (1) **First**, we will aim to learn what is referred to as the **(Stein) score function** of the log data density.
- (2) **Second**, Given this score function, we will be able to draw samples from $p_\theta(x)$ using Langevin Monte Carlo.

Going through these two steps...

Score-based Generative Models

First, we will aim to learn the **(Stein) score function** of the log data density, *i.e.*,

$$s_\theta(x) = \nabla_x \log p_\theta(x)$$

Score-based Generative Models

First, we will aim to learn the **(Stein) score function** of the log data density, *i.e.*,

$$s_\theta(x) = \nabla_x \log p_\theta(x)$$

Note that:

- The term *score* is also sometimes used to refer to the gradient of the log-likelihood function with respect to the parameter θ .
- We've seen this term before! (Recall in LMC/HMC...).

Score-based Generative Models

First, we will aim to learn the **(Stein) score function** of the log data density, i.e.,

$$s_\theta(x) = \nabla_x \log p_\theta(x)$$

This score function can also be written:

$$\begin{aligned} s_\theta(x) &= \nabla_x \log p_\theta(x) = \nabla_x \log \left(\frac{1}{Z(\theta)} \tilde{p}_\theta(x) \right) \\ &= -\nabla_x f_\theta(x) - \underbrace{\nabla_x \log Z(\theta)}_{=0} = -\nabla_x f_\theta(x) \end{aligned}$$

Score-based Generative Models

First, we will aim to learn the **(Stein) score function** of the log data density, *i.e.*,

$$s_\theta(x) = \nabla_x \log p_\theta(x)$$

This score function can also be written:

$$\begin{aligned} s_\theta(x) &= \nabla_x \log p_\theta(x) = \nabla_x \log \left(\frac{1}{Z(\theta)} \tilde{p}_\theta(x) \right) \\ &= -\nabla_x f_\theta(x) - \underbrace{\nabla_x \log Z(\theta)}_{=0} = -\nabla_x f_\theta(x) \end{aligned}$$

Definition of
unnormalized
PDF

Score-based Generative Models

First, we will aim to learn the **(Stein) score function** of the log data density, i.e.,

$$s_\theta(x) = \nabla_x \log p_\theta(x)$$

This score function can also be written:

$$\begin{aligned} s_\theta(x) &= \nabla_x \log p_\theta(x) = \nabla_x \log \left(\frac{1}{Z(\theta)} \tilde{p}_\theta(x) \right) \\ &= -\nabla_x f_\theta(x) - \underbrace{\nabla_x \log Z(\theta)}_{=0} = -\nabla_x f_\theta(x) \end{aligned}$$

Definition of
energy-based model

Definition of
unnormalized
PDF

Score-based Generative Models

First, we will aim to learn the **(Stein) score function** of the log data density, i.e.,

$$s_\theta(x) = \nabla_x \log p_\theta(x)$$

This score function can also be written:

$$\begin{aligned} s_\theta(x) &= \nabla_x \log p_\theta(x) = \nabla_x \log \left(\frac{1}{Z(\theta)} \tilde{p}_\theta(x) \right) \\ &= -\nabla_x f_\theta(x) - \underbrace{\nabla_x \log Z(\theta)}_{=0} = -\nabla_x f_\theta(x) \end{aligned}$$

Definition of
energy-based model

Definition of
unnormalized
PDF

Score-based Generative Models

Suppose we have learned an approximation $s_\theta(x)$ to the true score function:

$$s_\theta(x) \approx \nabla_x \log p(x) = \nabla_x \log \tilde{p}(x)$$

*I.e., it approximates the score function of
the distribution that generated our data*

Score-based Generative Models

Suppose we have learned an approximation $s_\theta(x)$ to the true score function:

$$s_\theta(x) \approx \nabla_x \log p(x) = \nabla_x \log \tilde{p}(x)$$

I.e., it approximates the score function of the distribution that generated our data

Recall that in Langevin Monte Carlo (SMC):

- In order to sample from a PDF, you just need access to the gradient of the unnormalized PDF!

Score-based Generative Models

Quick LMC Recap:

Gradient Based MCMC

Suppose we have a given PDF $p(x)$, $x \in \mathbb{R}^n$.

Consider an objective function: $f(x) = -\log p(x)$.

The gradient of this objective function is:

$$\nabla_x f(x) = \left[\frac{d}{dx_1} f(x), \dots, \frac{d}{dx_n} f(x) \right]^\top$$

Score-based Generative Models

Quick LMC Recap:

Gradient Based MCMC – Langevin Monte Carlo (LMC)

Langevin Monte Carlo (LMC) is very similar!

Each step in LMC can be written:

$$x_t = x_{t-1} - \frac{\epsilon_t}{2} \nabla_x f(x_{t-1}) + \sqrt{\epsilon_t} \mathcal{N}(0, I_n)$$

$$x_t = x_{t-1} - \eta_t \nabla_x f(x_{t-1})$$

Compare with gradient descent



Main difference is an extra
“stochastic term”

Can view this as “gradient
descent plus noise”

Score-based Generative Models

Therefore (after learning the score function):

Score-based Generative Models

Therefore (after learning the score function):

Second, we can simply run LMC with our learned score function $s_\theta(x)$:

$$x_{i+1} \leftarrow x_i + \epsilon \nabla_x s_\theta(x) + \sqrt{2\epsilon} \mathcal{N}(0, I), \quad i = 0, 1, \dots, K,$$

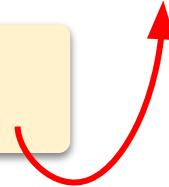
Score-based Generative Models

Therefore (after learning the score function):

Second, we can simply run LMC with our learned score function $s_\theta(x)$:

$$x_{i+1} \leftarrow x_i + \epsilon \nabla_x s_\theta(x) + \sqrt{2\epsilon} \mathcal{N}(0, I), \quad i = 0, 1, \dots, K,$$

Note the learned score
function is swapped in.



Score-based Generative Models

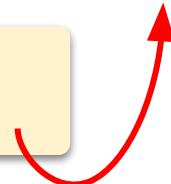
Therefore (after learning the score function):

Second, we can simply run LMC with our learned score function $s_\theta(x)$:

$$x_{i+1} \leftarrow x_i + \epsilon \nabla_x s_\theta(x) + \sqrt{2\epsilon} \mathcal{N}(0, I), \quad i = 0, 1, \dots, K,$$

Note the learned score
function is swapped in.

This is equivalent to previous
LMC, but step size is written
slightly different here.



Score-based Generative Models – An illustration

An illustration:

Score-based Generative Models – An illustration

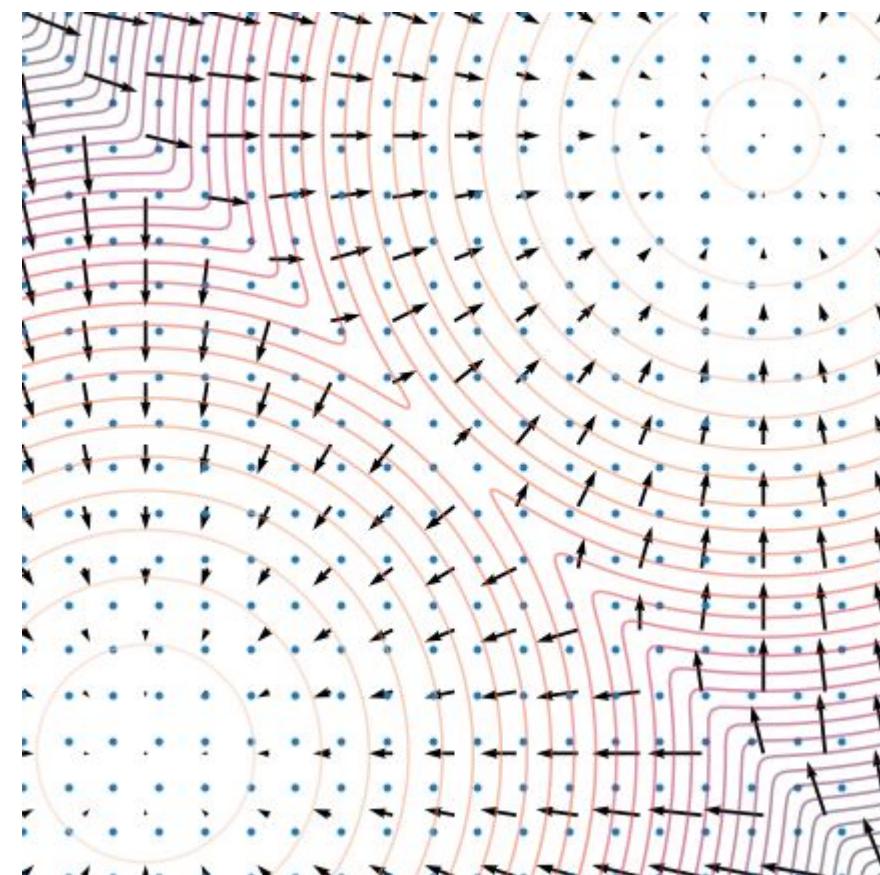
An illustration:

Contour lines denote the PDF $p_\theta(x)$.

Black arrows denote the score function $s_\theta(x)$.

Blue dots denote the samples produced via LMC.

(After random initialization).



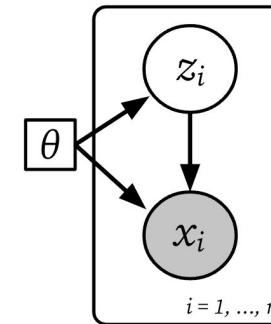
Score-based Generative Models – Intuition on Learning

What are we learning here (e.g., with a neural network) versus in VAEs?

Score-based Generative Models – Intuition on Learning

What are we learning here (e.g., with a neural network) versus in VAEs?

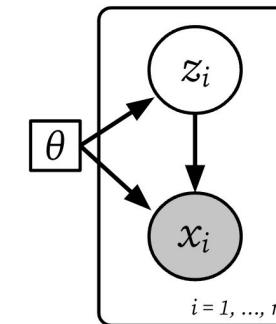
- In VAEs, we aim to learn a **direct map** from noise to data.
 - *I.e.*, directly learn the PDF of interest.
 - \Rightarrow a neural network that generates a sample from this PDF.



Score-based Generative Models – Intuition on Learning

What are we learning here (e.g., with a neural network) versus in VAEs?

- In VAEs, we aim to learn a **direct map** from noise to data.
 - I.e., directly learn the PDF of interest.
 - \Rightarrow a neural network that generates a sample from this PDF.
- In score-based models, we aim to learn a **score function**
 - Which we can view as pointing us in the direction of steepest ascent of PDF.
 \Rightarrow Use neural network *repeatedly* to take gradient steps to sample from this PDF.



Score-based Generative Models

So the main question is: How do we learn an approximation $s_\phi(x)$ to the score function $s_\theta(x)$?

Score-based Generative Models

So the main question is: How do we learn an approximation $s_\phi(x)$ to the score function $s_\theta(x)$?

⇒ Answer: we will minimize the (so called) *Fisher divergence* between the model and data distributions, written as:

$$\mathcal{L}(\theta) = \mathbb{E}_{p(x)} [\|\nabla_x \log p(x) - s_\theta(x)\|_2^2]$$

Intuitively: this compares the squared L2 distance between the ground-truth data score and the score-based model.

Score-based Generative Models

So the main question is: How do we learn an approximation $s_\phi(x)$ to the score function $s_\theta(x)$?

⇒ Answer: we will minimize the (so called) *Fisher divergence* between the model and data distributions, written as:

$$\mathcal{L}(\theta) = \mathbb{E}_{p(x)} [\|\nabla_x \log p(x) - s_\theta(x)\|_2^2]$$

However, there is an issue: we don't know the ground truth PDF $p(x)$ and ground truth data score $\nabla_x \log p(x)$.

Score-based Generative Models

Solution: Use a technique called **Score Matching**.

This will let us minimize the Fisher divergence without knowledge of the ground-truth data score.

Journal of Machine Learning Research 6 (2005) 695–709
Submitted 11/04; Revised 3/05; Published 4/05

Estimation of Non-Normalized Statistical Models by Score Matching

Aapo Hyvärinen

Helsinki Institute for Information Technology (BRU)
Department of Computer Science
FIN-00014 University of Helsinki, Finland

AAPO.HYVARINEN@HELSINKI.FI

Editor: Peter Dayan

Abstract

One often wants to estimate statistical models where the probability density function is known only up to a multiplicative normalization constant. Typically, one then has to resort to Markov Chain Monte Carlo methods, or approximations of the normalization constant. Here, we propose that such models can be estimated by minimizing the expected squared distance between the gradient of the log-density given by the model and the gradient of the log-density of the observed data. While the estimation of the gradient of log-density function is, in principle, a very difficult non-parametric problem, we prove a surprising result that gives a simple formula for this objective function. The density function of the observed data does not appear in this formula, which simplifies to a sample average of a sum of some derivatives of the log-density given by the model. The validity of the method is demonstrated on multivariate Gaussian and independent component analysis models, and by estimating an overcomplete filter set for natural image data.

Keywords: statistical estimation, non-normalized densities, pseudo-likelihood, Markov chain Monte Carlo, contrastive divergence

1. Introduction

In many cases, probabilistic models in machine learning, statistics, or signal processing are given in the form of non-normalized probability densities. That is, the model contains an unknown normalization constant whose computation is too difficult for practical purposes.

Assume we observe a random vector $\mathbf{x} \in \mathbb{R}^n$ which has a probability density function (pdf) denoted by $p_{\mathbf{x}}(\cdot)$. We have a parametrized density model $p(\cdot; \boldsymbol{\theta})$, where $\boldsymbol{\theta}$ is an m -dimensional vector of parameters. We want to estimate the parameter $\boldsymbol{\theta}$ from \mathbf{x} , i.e. we want to approximate $p_{\mathbf{x}}(\cdot)$ by $p(\cdot; \hat{\boldsymbol{\theta}})$ for the estimated parameter value $\hat{\boldsymbol{\theta}}$. (We shall here consider the case of continuous-valued variables only.)

The problem we consider here is that we only are able to compute the pdf given by the model up to a multiplicative constant $Z(\boldsymbol{\theta})$:

$$p(\mathbf{x}; \boldsymbol{\theta}) = \frac{1}{Z(\boldsymbol{\theta})} q(\mathbf{x}; \boldsymbol{\theta}).$$

That is, we do know the functional form of q as an analytical expression (or any form that can be easily computed), but we do *not* know how to easily compute Z which is given by

©2005 Aapo Hyvärinen.

Next Class

Next Class

Next class, we will cover:

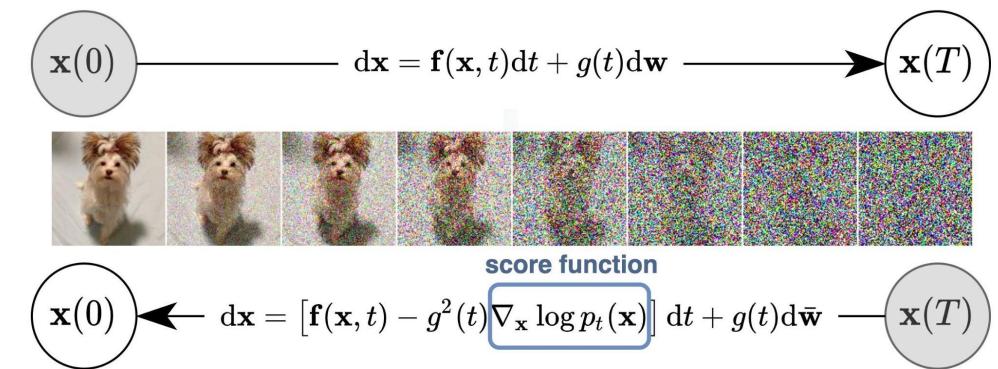
Next Class

Lecture: (Continued) score-based generative models and diffusion models.

Next Class

Lecture: (Continued) score-based generative models and diffusion models.

- The (Stein) score function, score matching.
- Score-based generative models.
- Sampling with Langevin Monte Carlo.
- Denoising diffusion probabilistic models.
- Learning with variational inference.

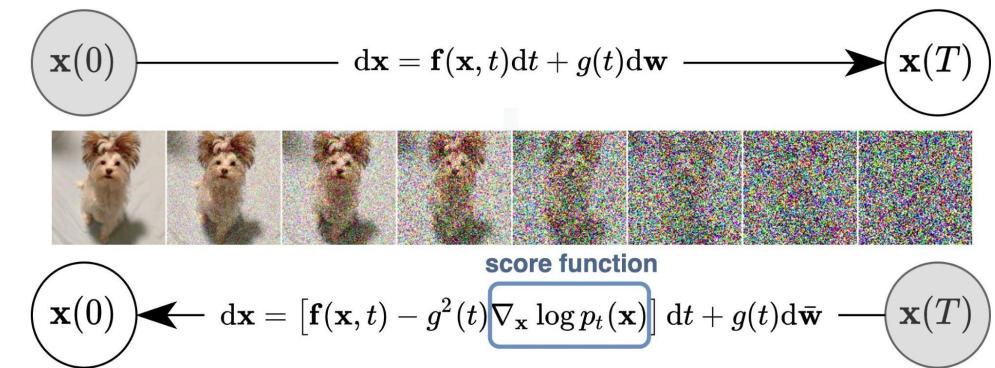


Source: Yang Song, "Generative Modeling by Estimating Gradients of the Data Distribution"

Next Class

Lecture: (Continued) score-based generative models and diffusion models.

- The (Stein) score function, score matching.
- Score-based generative models.
- Sampling with Langevin Monte Carlo.
- Denoising diffusion probabilistic models.
- Learning with variational inference.



Source: Yang Song, "Generative Modeling by Estimating Gradients of the Data Distribution"

After:

- Paper presentations from students.

