

CSCI 699 - ProbGen

Probabilistic and Generative Models

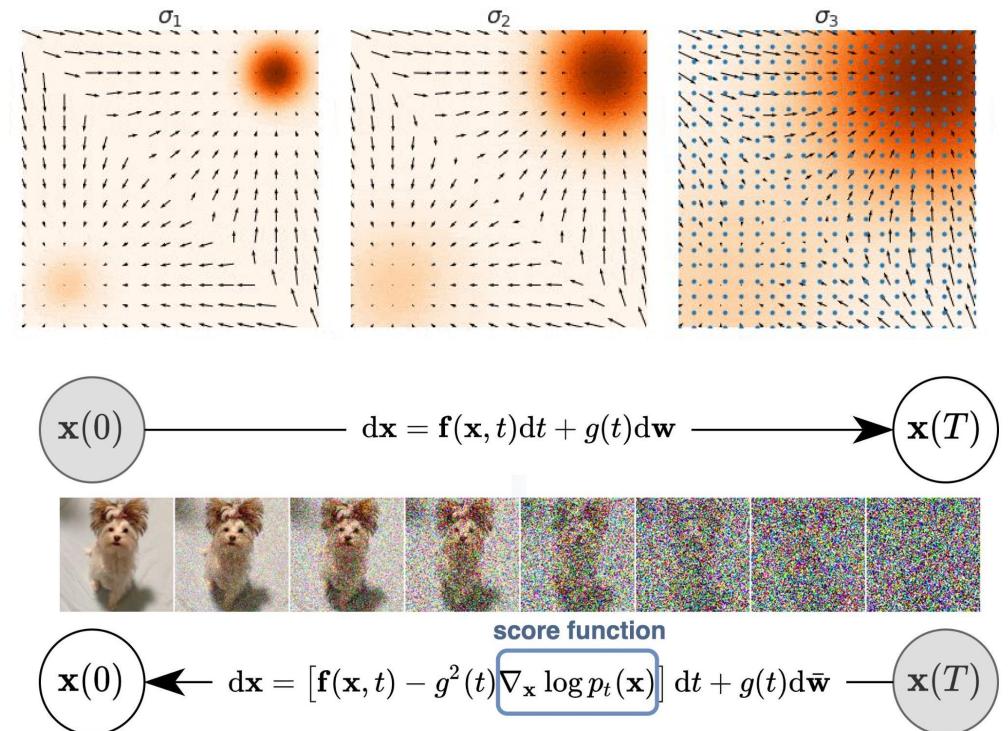
Willie Neiswanger

Lecture 8 - Score-based and Diffusion Generative Models

Today

Today

Lecture: Score-based generative models, and intro to diffusion models.

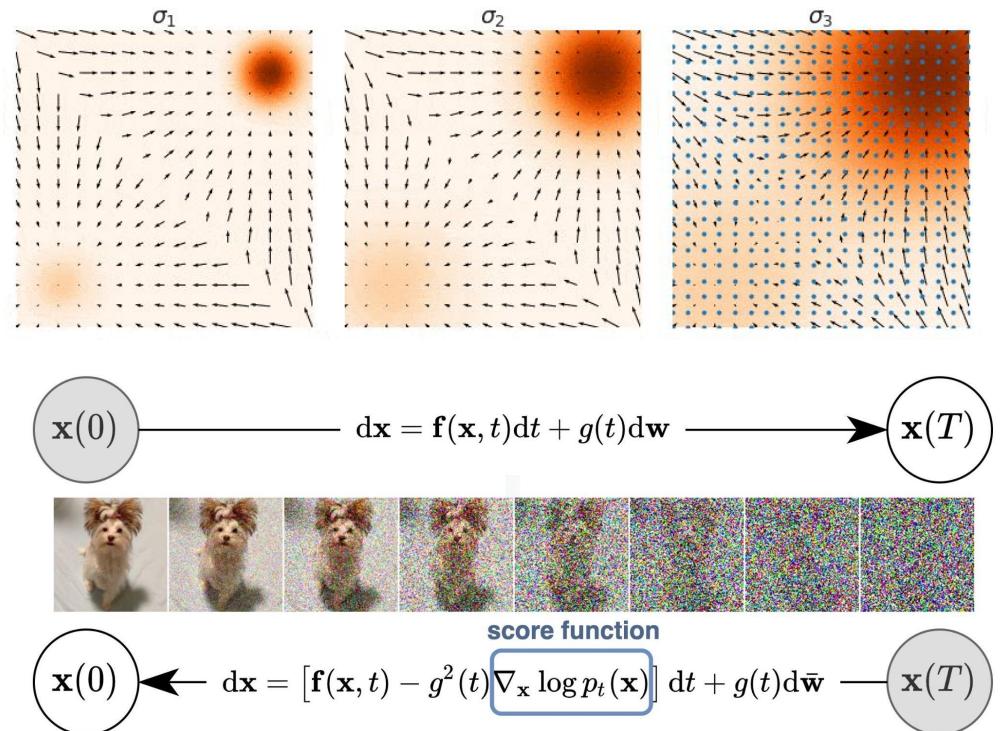


Source: Yang Song, "Generative Modeling by Estimating Gradients of the Data Distribution"

Today

Lecture: Score-based generative models, and intro to diffusion models.

- Score-based Generative Models
 - Modeling the (Stein) Score Function
 - Annealed Langevin Monte Carlo

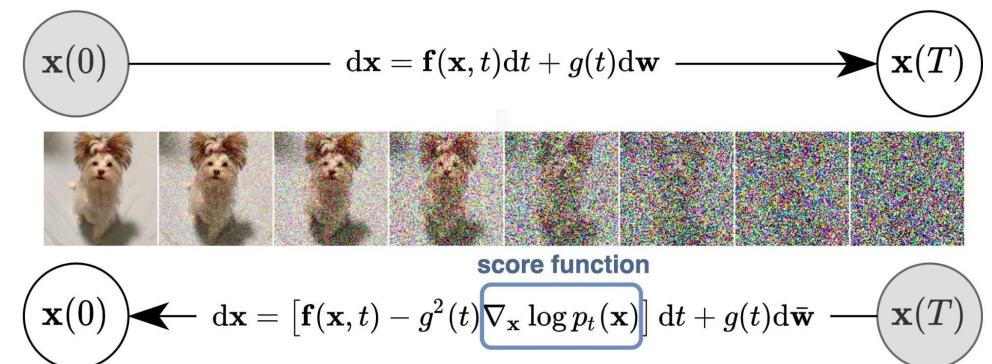
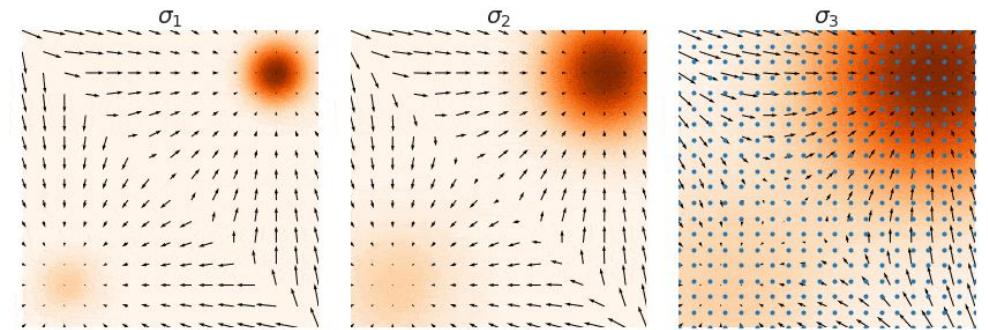


Source: Yang Song, "Generative Modeling by Estimating Gradients of the Data Distribution"

Today

Lecture: Score-based generative models, and intro to diffusion models.

- Score-based Generative Models
 - Modeling the (Stein) Score Function
 - Annealed Langevin Monte Carlo
- Framing as a diffusion model.
- VI for diffusion models.

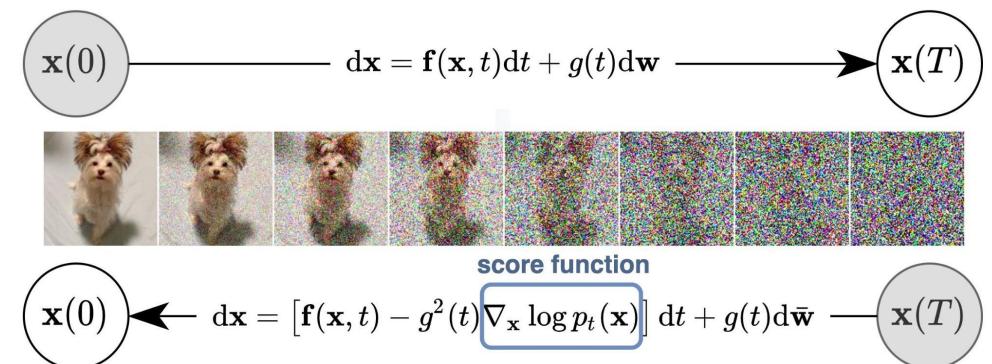
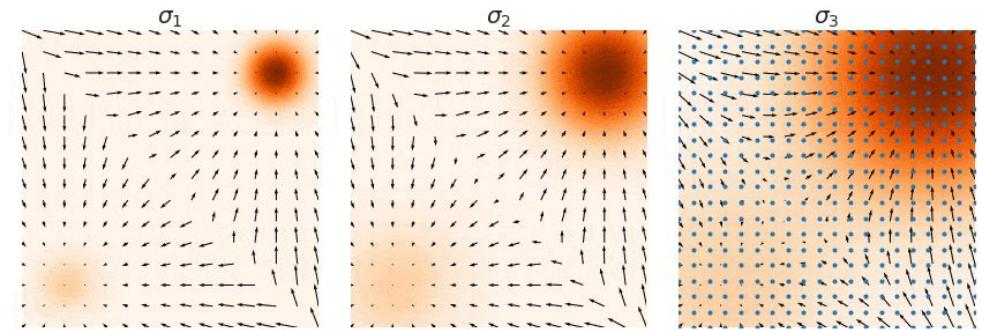


Source: Yang Song, "Generative Modeling by Estimating Gradients of the Data Distribution"

Today

Lecture: Score-based generative models, and intro to diffusion models.

- Score-based Generative Models
 - Modeling the (Stein) Score Function
 - Annealed Langevin Monte Carlo
- Framing as a diffusion model.
- VI for diffusion models.



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

After:

- Second batch of four paper presentations.

Today

After: Second batch of four paper presentations.

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

arXiv:2005.00341v1 [eess.AS] 30 Apr 2020

arXiv:1505.05700v6 [stat.ML] 14 Jun 2016

arXiv:2501.19393v3 [cs.CL] 1 Mar 2025

Today

After: Second batch of four paper presentations.

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

arXiv:2005.00341v1 [eess.AS] 30 Apr 2020

arXiv:1505.05706v6 [stat.ML] 14 Jun 2016

s1: Simple test-time scaling

Average thinking time (tokens)	MATH500 (%)	AIME24 (%)	GRQA (%)
512	~85	~20	~40
1024	~95	~60	~60
2048	~95	~60	~60
4096	~95	~60	~60

Figure 1. Test-time scaling with s1-32B. We benchmark s1-32B on reasoning-intensive tasks and very fast test-time compute.

Test-time scaling is a promising new approach to language modeling that uses extra test-time compute to improve performance. Recently, OpenAI's s1 model showed this capability but did not publicly share its methodology. In many replication efforts we have checked the implementation to achieve test-time scaling and strong reasoning performance. First, we curate a small dataset s1KB of 1,000 questions paired with reasoning traces relying on GEMINI. Second, we split them into three dimensions: difficulty, diversity, and quality. Second, we develop budget forcing to control test-time compute by forcefully terminating the model's thinking process or lengthening it by appending "wait" multiple times until the budget is used up. When it times out, the user can lead the model to double-check its answer, often fixing incorrect reasoning steps. After supervised finetuning the Gemini-2.5-32B large language model on s1KB and equipping it with budget forcing, we find that the model can execute tests of reasoning on competition math questions by up to 27% (MATH and AIME24). Further, scaling s1-32B with budget forcing allows extrapolating beyond the current test-time compute constraint from 50% to 57% on AIME24. Our model, data, and code are open-source at <https://github.com/simple-learning/s1>.

1. Introduction

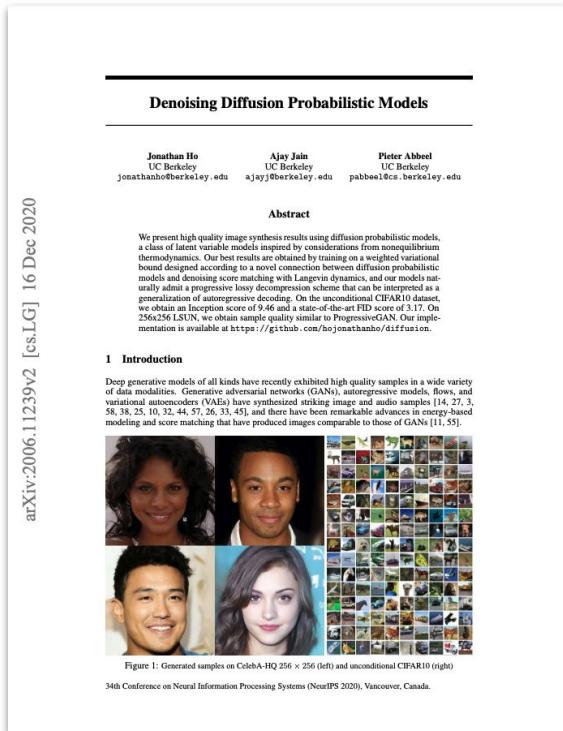
Performance improvements of language models (LMs) over the past years have largely relied on scaling up train-time compute using large-scale self-supervised pretraining (Kaggle et al., 2020; Hoffmann et al., 2022). The creation of these powerful models has set the stage for a new scaling paradigm built on top of their *test-time scaling*. The aim

of this approach is to increase the compute at test time to get better results. There has been much work exploring this idea (Snoek et al., 2024; Welelck et al., 2024), and the viability of this paradigm was recently validated by OpenAI (OpenAI, 2024). It has been shown that scaling reasoning performance with comment grants scales from test-time compute. OpenAI describes their approach as using large-scale reinforcement learning (RL) implying the use of sizable amounts of data (OpenAI, 2024). This has led to various autoregressive approaches such as beam search-like Monte Carlo Tree Search (Guo et al., 2024b; Zhang et al., 2024a), multi-agent approaches (Qin et al., 2024a), and others (Wang et al., 2024c; Huang et al., 2024b; 2025). Among these approaches, the one proposed by OpenAI et al. (2024) has successfully replicated o1-level performance by employing reinforcement learning via millions of samples and multiple training stages. However, despite the large number of o1 replication attempts, none have openly replicated a test-time scaling model on a public dataset. Thus, we ask: what is the simplest approach to achieve both test-time scaling and strong reasoning performance?

We show that training on only 1,000 samples with next-token prediction and controlling thinking duration via a simple test-time technique we refer to as *budget forcing* leads to a state-of-the-art reasoning model that scales in performance with more time compute. Specifically, we compare s1, which consists of 1,000 carefully curated questions paired with reasoning traces and answers distilled from Gemini Thinking Experimental (Google, 2024). We perform supervised fine-tuning (SFT) of an off-the-shelf pretrained model

Today

After: Second batch of four paper presentations.



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

Denoising Diffusion Probabilistic Models

Jonathan Ho
UC Berkeley
jonathanho@berkeley.edu Ajay Jain
UC Berkeley
ajay@berkeley.edu Pieter Abbeel
UC Berkeley
pabbeel@cs.berkeley.edu

Abstract

We present high-quality images synthesized using denoising probabilistic models, a class of latent variable models inspired by considerations from nonequilibrium thermodynamics. Our best results are obtained by training on a weighted variational bound designed according to a novel connection between diffusion probabilistic models and generative models trained with Langevin dynamics. Our models naturally admit a progressive lossy decoding scheme that can be interpreted as a generalization of autoregressive decoding. On the unconditional CIFAR10 dataset, we obtain an Inception score of 9.46 and a state-of-the-art FID score of 3.17. On 256x256 LSUN, we obtain sample quality similar to ProgressiveGAN. Our implementation is available at <https://github.com/jonathanho/diffusion>.

1 Introduction

Deep generative models of all kinds have recently exhibited high quality samples in a wide variety of applications. Generative adversarial networks (GANs), autoregressive models, flows, and variational autoencoders (VAEs) have synthesized images and audio samples [14, 27, 3, 58, 32, 25, 10, 32, 44, 57, 26, 33, 45], where there have been remarkable advances in encoding modeling and score matching that have produced images comparable to those of GANs [11, 35].



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.



1. Introduction

Music is an integral part of human culture, existing from the earliest periods of human civilization and evolving into a wide diversity of forms. It is the expression of whether the human spirit in its evolution, the expression of whether the human spirit can capture this creative process has fascinated computer scientists for decades. We have had algorithms generating piano sheet music (Hiller Jr & Isaacson, 1957; Moorer, 1972; Hadjeres et al., 2017; Huang et al., 2018), digital vocoders (Boulanger-Lewandowski & Schmidhuber, 2001; Engels et al., 2006; Blauw & Bonada, 2017) and also synthesizers producing timbres for various musical instruments (Engel et al., 2017; 2019). Each captures a specific aspect of music generation: melody, composition, timbre, and the human voice singing. However, a single system to do it all remains elusive.

The field of generative models has made tremendous progress in the last few years. One of the aims of generative modeling is to capture the salient aspects of the data and to generate new samples that are indistinguishable from the true. The hypothesis is that by learning to generate data we can learn the best features of the data. We are surrounded by highly complex distributions in the visual, audio, and text domain, and in recent years we have developed

In this work, we show that we can use state-of-the-art deep generative models to produce a single system capable of generating diverse high-fidelity music in the raw audio domain, with long-range coherence spanning multiple minutes. Our approach uses a hierarchical VQ-VAE architecture (Razavi et al., 2019).

Figure 2: Jukebox: A Generative Model for Music. This figure includes the arXiv ID (arXiv:2005.00341v1 [stat.ML] 14 Jun 2016), the title 'Jukebox: A Generative Model for Music', the authors (Danilo Jimenez Rezende, Shakir Mohamed), and the abstract.

Abstract

The choice of approximate posterior distribution is one of the core problems in variational inference. Most applications of variational inference employ mean-field or simple structured approximations in order to allow for efficient inference, focusing on mean-field or other simple structured approximations. The choice of approximation has a significant impact on the quality of inference and the quality of the generated samples. We introduce a new approach for specifying flexible, arbitrarily complex and scalable approximate posterior distributions. Our approximations are distributions constructed through a normative flow, where a simple initial distribution is transformed into a more complex one by applying a sequence of invertible transformations until a desired level of complexity is attained. This makes the problem more challenging, as the space of audio is extremely high dimensional with a high amount of information content to model. There has been much success, while developing approximate posteriors in the audio domain (Oord et al., 2016; Mehrilahi et al., 2017; Yamamoto et al., 2020) or in the spectrogram domain (Vasquez & Lewis, 2019). The key bottleneck is that modeling the raw audio directly introduces extremely long-range dependencies, making it computationally challenging to learn a high-dimensional space. As a way to reduce the difficulty is to learn a lower-dimensional encoding of the audio with the goal of losing the less important information while still maintaining the most of the musical information. This approach has demonstrated some success in generating short instrumental pieces produced to a set of a few instruments (Oord et al., 2017; Dieleman et al., 2018).

In this work, we show that we can use state-of-the-art deep generative models to produce a single system capable of generating diverse high-fidelity music in the raw audio domain, with long-range coherence spanning multiple minutes. Our approach uses a hierarchical VQ-VAE architecture (Razavi et al., 2019).

^{*}Richard Feynmann famously said, "What I cannot create, I do not understand."

[†]Equal contribution. [†]OpenAI, San Francisco. Correspondence to: cjukebox@openai.com.



1. Introduction

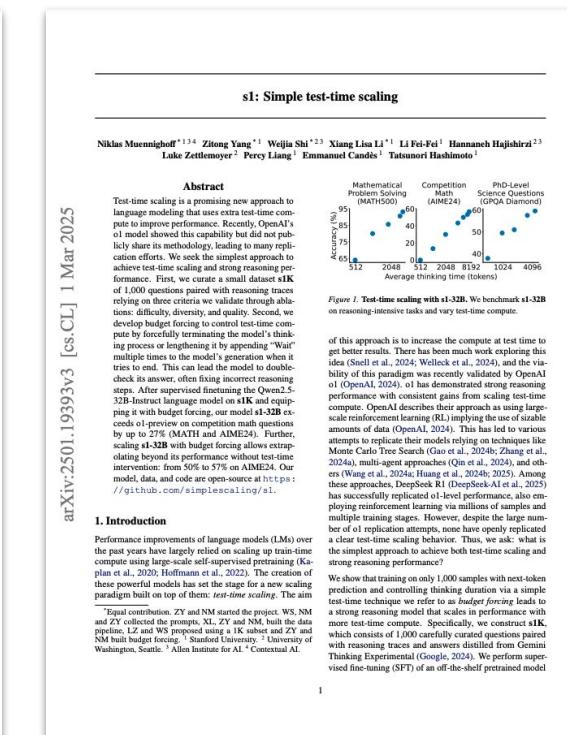
There has been a great deal of renewed interest in variational inference as a means of scaling probabilistic modeling to increasingly complex problems on increasingly large datasets. Variational inference has been at the core of large-scale topic models of text (Hoffmann et al., 2013), provides the state-of-the-art in semi-supervised classification (Kingma et al., 2014), drives the models that currently produce the most robust generative models of images (Gregor et al., 2014; Kingma & Welling, 2013; Kingma & Welling, 2014), and is a default tool for the understanding

of many physical and chemical systems. Despite these successes and ongoing advances, there are a number of disadvantages of variational methods that limit their power and handling. We here adopt a new method for statistical inference. It is of course the limitation, the choice of posterior approximation, that we address in this paper.

Variational inference requires that intractable posterior distributions be approximated by a class of known probability distributions over which we search for the best approximation. These posterior distributions are often intractable and often limited (e.g. mean-field approximations), implying that no solution is ever able to resemble the true posterior distribution. This is a widely raised objection to variational methods, in that we sacrifice inferential methods such as posterior sampling through a normalizing flow, where a single sample does not transform into a more complex one by applying a sequence of invertible transformations until a desired level of complexity is attained. This makes the problem more challenging, as the space of audio is extremely high dimensional with a high amount of information content to model. There has been much success, while developing approximate posteriors in the audio domain (Oord et al., 2016; Mehrilahi et al., 2017; Yamamoto et al., 2020) or in the spectrogram domain (Vasquez & Lewis, 2019). The key bottleneck is that modeling the raw audio directly introduces extremely long-range dependencies, making it computationally challenging to learn a high-dimensional space. As a way to reduce the difficulty is to learn a lower-dimensional encoding of the audio with the goal of losing the less important information while still maintaining the most of the musical information. This approach has demonstrated some success in generating short instrumental pieces produced to a set of a few instruments (Oord et al., 2017; Dieleman et al., 2018).

There is much evidence that richer, more faithful posterior approximations do result in better performance. For example, when compared to sigmoid belief networks that make use of mean-field approximations, deep autoregressive models with variational approximations with a more expressive dependency structure that provides a clear improvement in performance (Mnih & Gregor, 2014). There is also a large body of evidence that the more expressive a posterior approximation, the better the performance. For example, in the context of variational inference, Wing et al. (2011) provide an comparison of two commonly used approaches. The first is the widely-observed problem of under-estimation of the variance of the posterior distribution, which can lead to poor performance in a reasonable decision based on the chosen posterior approximation. The second is that the limited capacity of the posterior approximation can also result in biases in the MAP estimates of any model parameters (and this is the case e.g. in linear-separable classification).

A number of proposals for rich posterior approximations have been explored, typically based on structured mean-field approximations that incorporate some basic form of dependency within the approximate posterior. Another potentially powerful alternative would be to specify an approximate posterior as a mixture model, such as those developed by Jaakkola & Jordan (1998); Jordan et al. (1999); Gerstenman et al. (2012). But the mixture approach limits



Today

After: Second batch of four paper presentations.

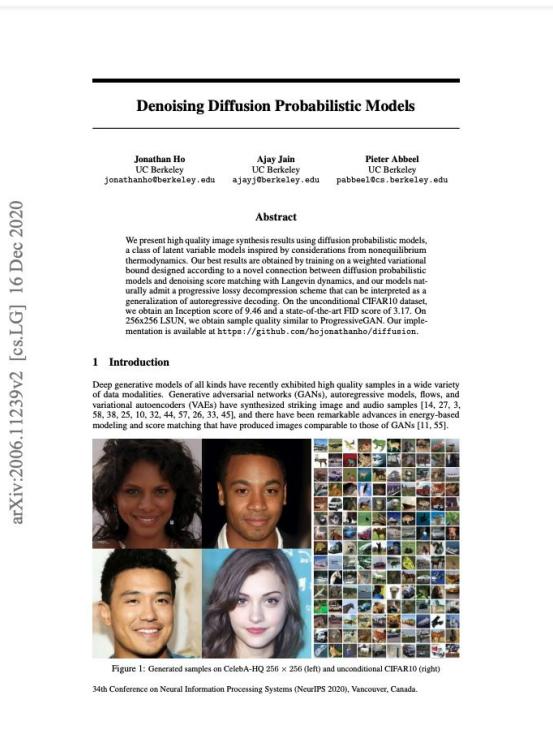


Figure 1: Generated samples on CelebA-HQ 256 × 256 (left) and unconditional CIFAR10 (right)

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

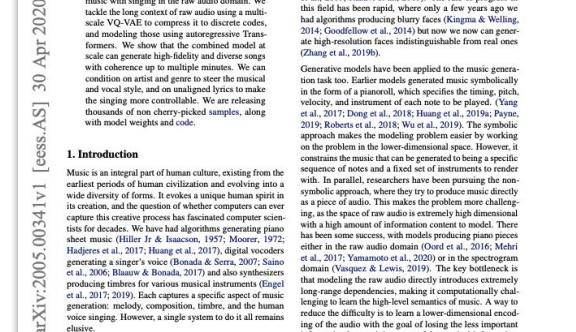
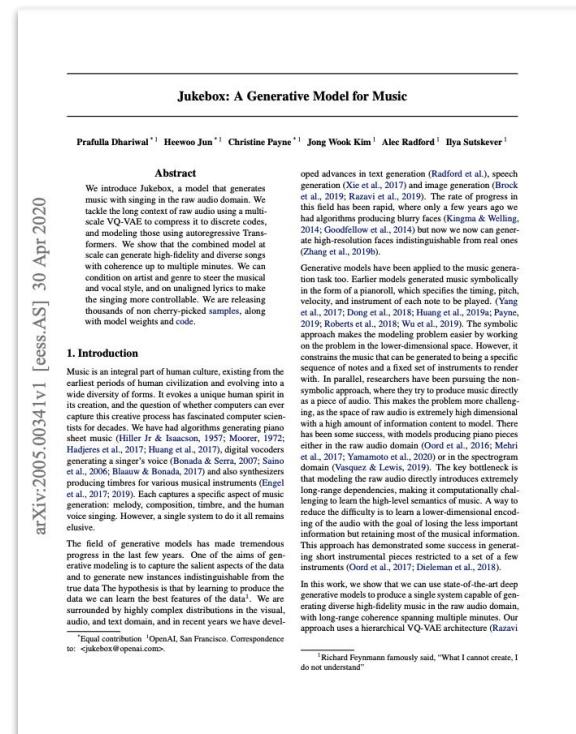
We present high-quality images synthesized using denoising probabilistic models, a class of latent variable models inspired by considerations from nonequilibrium thermodynamics. Our best results are obtained by training on a weighted variational bound designed according to a novel connection between diffusion probabilistic models and generative models trained with Langevin dynamics. Our models naturally admit a progressive lossy decoding scheme that can be interpreted as a generalization of autoregressive decoding. On the unconditional CIFAR10 dataset, we obtain an Inception score of 9.46 and a state-of-the-art FID score of 3.17. On 256x256 LSUN, we obtain sample quality similar to ProgressiveGAN. Our implementation is available at <https://github.com/jonathanho/diffusion>.

1 Introduction

Deep generative models of all kinds have recently exhibited high quality samples in a wide variety of applications. Generative adversarial networks (GANs), autoregressive models, flows, and variational autoencoders (VAEs) have synthesized images and audio samples [14, 27, 3, 58, 32, 25, 10, 32, 44, 57, 26, 33, 45], while there have been remarkable advances in energy-based modeling and score matching that have produced images comparable to those of GANs [11, 35].



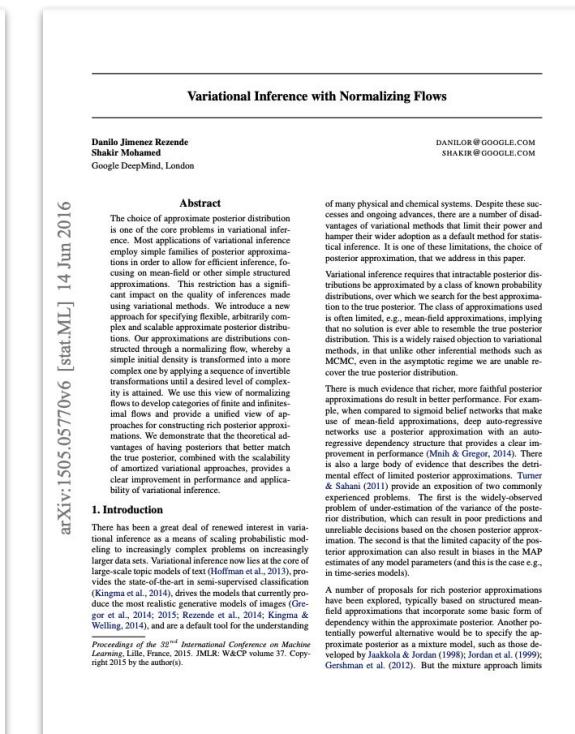
34th Conference on Neural Information Processing Systems (NeurIPS 2020), Vancouver, Canada.



¹Equal contribution. ¹OpenAI, San Francisco. Correspondence to: cjukebox@openai.com.

¹Richard Feynmann famously said, "What I cannot create, I do not understand."

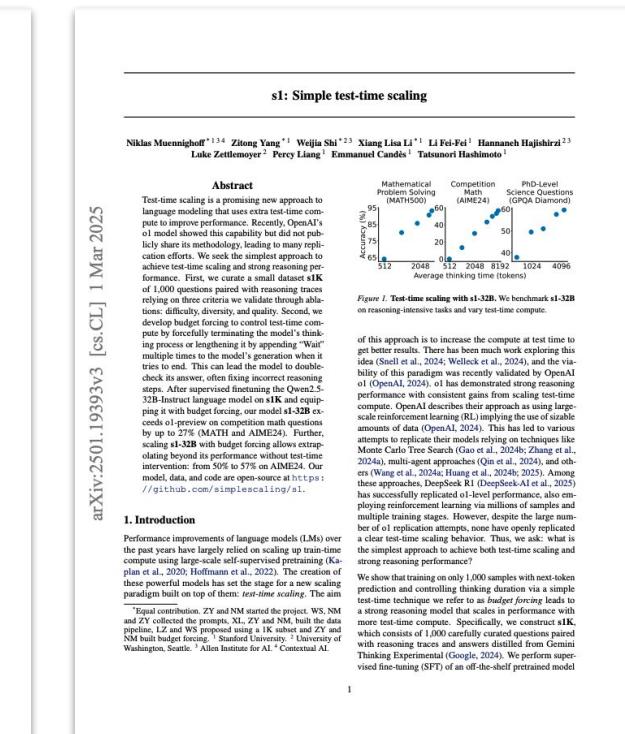
¹Xiang Li, Li Fei-Fei, Hannaneh Hajishirzi, and Percy Liang. 2015. JMLR: W&CP volume 37. Copyright 2015 by the author(s).



¹Equal contribution. ¹OpenAI, San Francisco. Correspondence to: cjukebox@openai.com.

¹Richard Feynmann famously said, "What I cannot create, I do not understand."

¹Xiang Li, Li Fei-Fei, Hannaneh Hajishirzi, and Percy Liang. 2015. JMLR: W&CP volume 37. Copyright 2015 by the author(s).



arXiv:2005.00341v1 [stat.ML] 14 Jun 2016

arXiv:1505.05770v6 [stat.ML] 30 Apr 2020

arXiv:1939v3 [cs.CL] 1 Mar 2025

s1: Simple test-time scaling

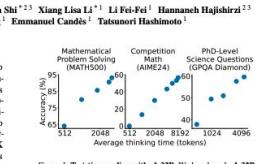


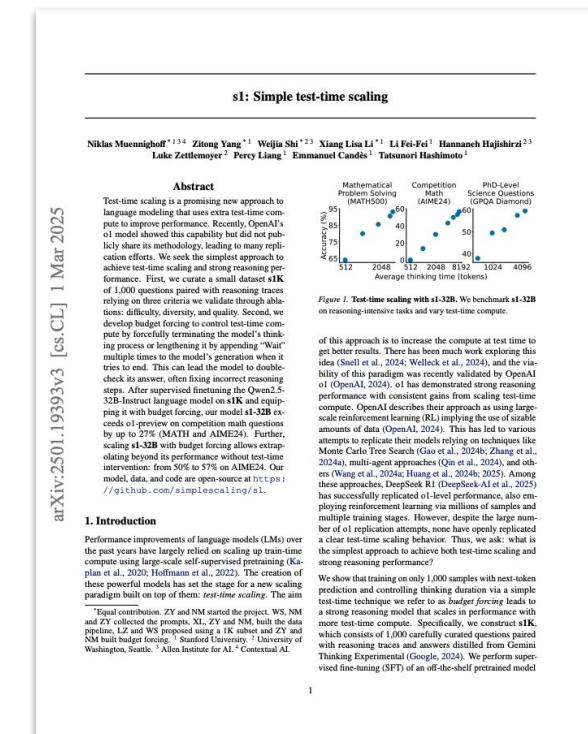
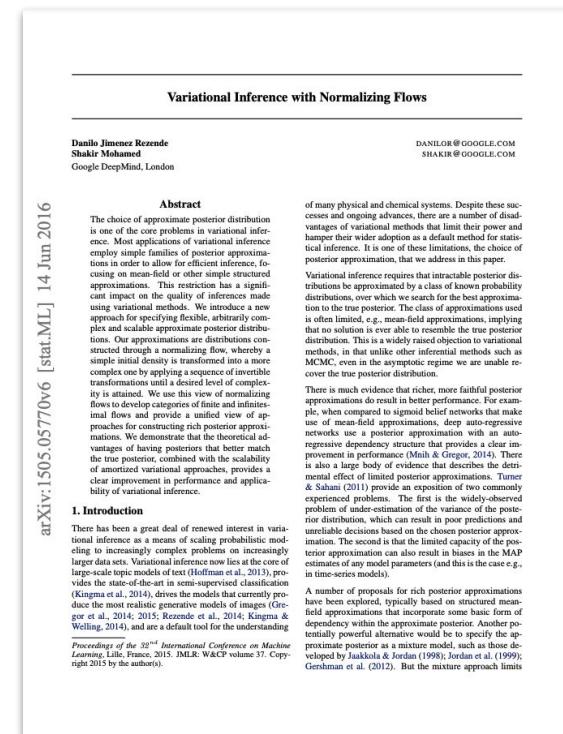
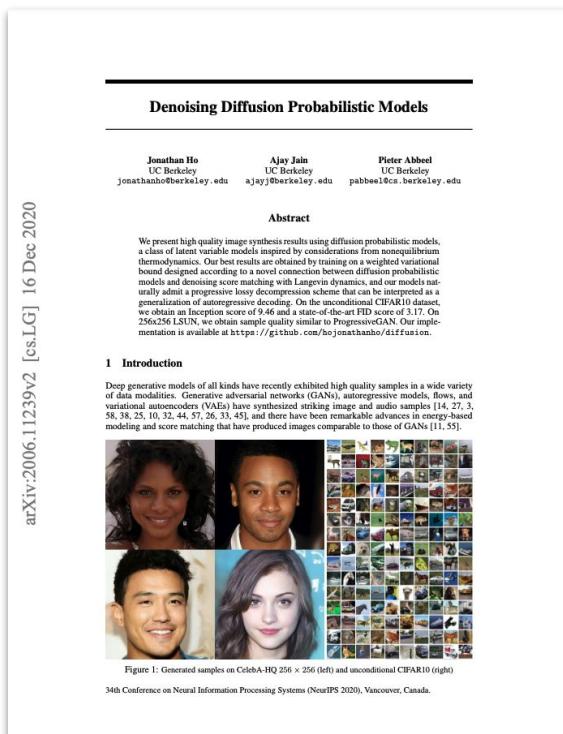
Figure 1: Test-time scaling with s1-32B. We benchmark s1-32B on reasoning-intensive tasks and vary test-time compute.

of this approach is to increase the compute at test time to generate results. There has been much work exploring this idea (Senni et al., 2024; Wolfschmidt et al., 2024), and the viability of this paradigm was recently validated by OpenAI's el model (OpenAI, 2024). el has demonstrated strong reasoning performance with consistent gains for scaling test-time compute. OpenAI demonstrated their approach as having large-scale reasoning gains by simply scaling the number of available amounts of data (OpenAI, 2024). This has led to various attempts to replicate their models relying on techniques like Monte Carlo Tree Search (Gao et al., 2024b; Zhang et al., 2024a), Reinforcement Learning (Qu et al., 2024), and Generators (Wang et al., 2024a; Huang et al., 2024b, 2025). Among these approaches, DeepSeek R1 (DeepSeek et al., 2025) has successfully replicated el-level performance, also employing reinforcement learning to learn the scaling of softmax masking training stages. However, despite a large number of el replication attempts, none have openly replicated a clear test-time scaling behavior. Thus, we ask: what is the simplest approach to achieve both test-time scaling and strong reasoning performance?

We propose *test-time scaling* via 1,000 samples with test-time pre-training and controlling training duration via a simple test-time technique we refer to as *budget forcing* leads to a strong reasoning model that scales in performance with more test-time compute. Specifically, we construct **s1K**, which contains 1,000 samples and can be paired with reasoning traces and answer distributions from our Thinking Experimental (Google, 2024). We perform supervised fine-tuning (SFT) of an off-the-shelf pre-trained model

Today

After: Second batch of four paper presentations.



s1: Simple test-time compute/scaling

Today

Quick note on scribe and discussion lead assignments:

- As you know, discussion leads bring ~5 questions on their assigned paper, and scribes write up ~1 page document.
- Both assignments should be uploaded as a **PDF file** on [Brightspace](#).
- Ideally within a ~week of your scribe/discussion lead duties.

Today

Here's how to upload:

Today

Here's how to upload: (1) Go to main course homepage on Brightspace.

The screenshot shows the main course homepage for "20251_30165 CSCI-699: Special To...". The page features a banner image of green leaves with water droplets. The navigation bar includes links for Home, Announcements, Content, Activities, My Grades, Help, Course Tools, and Library Resources. A user profile for "Willie Neiswanger_Test" is shown, along with impersonating and settings icons. The main content area contains several widgets:

- Slim Announcements Widget:** Displays two announcements:
 - "Lecture Slides Added and Note on Next Class" (Posted Wednesday, January 22, 2025 at 1:20 PM): Hi all, I've uploaded the slides from last Friday's class here on Brightspace (under Content > Lectures > Lecture 1 Slides). And a couple of reminders for this coming Friday's class: This week we will cover: An Introduction ... [Read More](#)
 - "Welcome to Probabilistic and Generative Models!" (Posted Thursday, January 16, 2025 at 1:14 PM): Welcome to CSCI 699: Probabilistic and Generative Models, Spring 2025! This is the private course materials and submission site. Here is the public website: <https://willieneis.github.io/probgen-spring2025/> [Show All Announcements](#)
- Calendar:** Shows "Friday, March 7, 2025" and "Upcoming events".
- Activity Feed:** Shows "Latest Posts": "There's nothing here just yet."
- Announcements:** Shows the same "Lecture Slides Added and Note on Next Class" announcement as the widget.
- Multi-Profile Widget:** Shows a placeholder for a user profile.

Today

Here's how to upload: (2) Click on *Activities > Assignments* in menu bar at top.

The screenshot shows the USC Viterbi School of Engineering Brightspace LMS homepage. At the top, the navigation bar includes links for Home, Announcements, Content, Activities (which is highlighted with a red box), My Grades, Help, Course Tools, and Library Resources. Below the navigation bar is a banner for "CSCI 699: Probabilistic and Generative Models". A dropdown menu for "Activities" is open, showing options for Assignments, Quizzes, Discussions, and Groups. The main content area contains several widgets: a "Slim Announcements Widget" with two entries about lecture slides and a welcome message; a "Calendar" widget showing "Friday, March 7, 2025" and "Upcoming events"; an "Activity Feed" which is currently empty; and a "Latest Posts" section which also contains no content. At the bottom, there is a "Multi-Profile Widget". The URL in the browser address bar is https://brightspace.usc.edu/courses/20251_30165/CSCI-699-Special-Topics-and-Research-Methods-in-Computer-Vision-and-Machine-Learning.

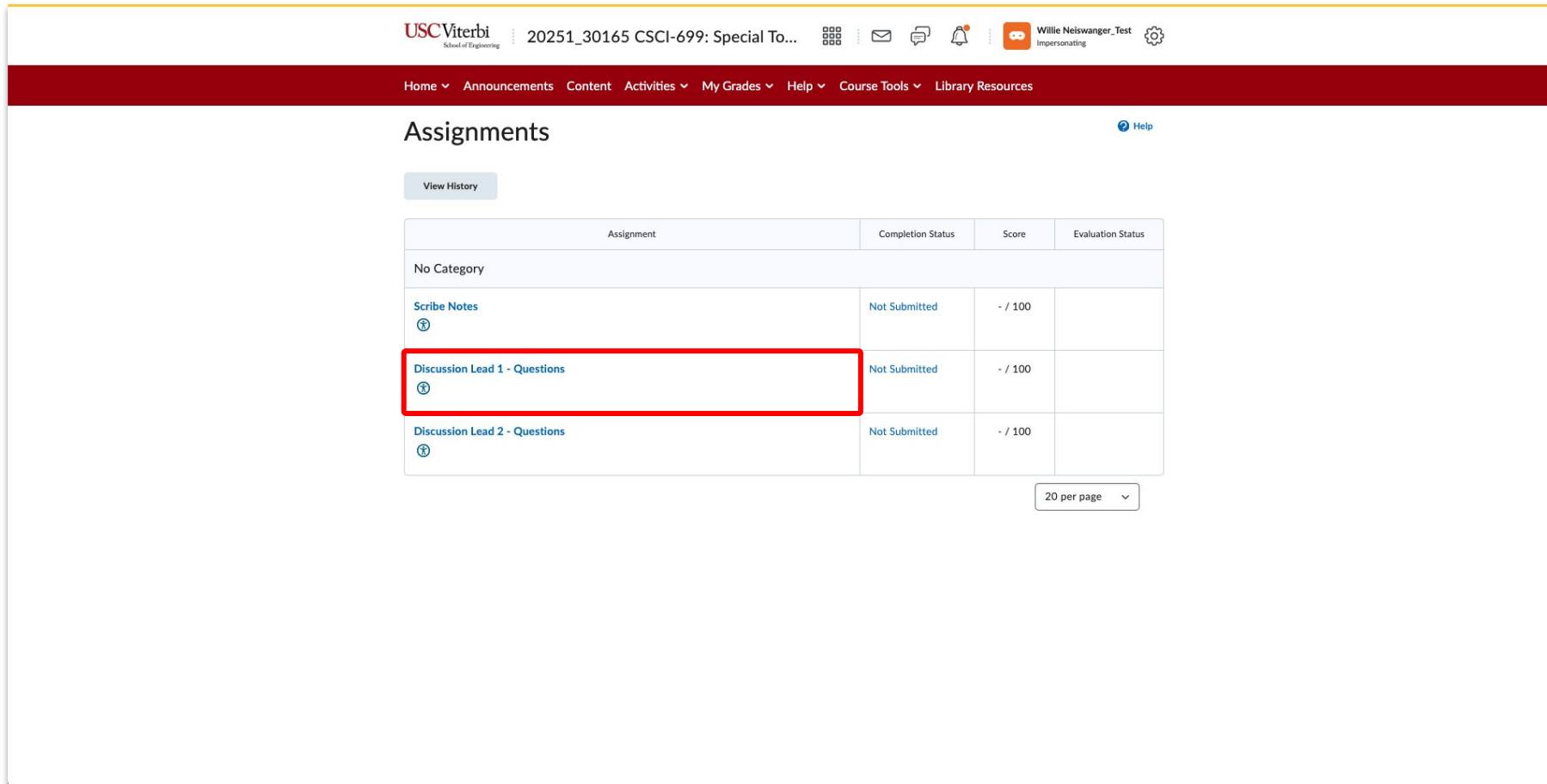
Today

Here's how to upload: (2) Click on *Activities > Assignments* in menu bar at top.

The screenshot shows the USC Viterbi School of Engineering Brightspace course page for CSCI 699: Probabilistic and Generative Models. The top navigation bar includes links for Home, Announcements, Content, Activities (with a dropdown menu), My Grades, Help, Course Tools, and Library Resources. The Activities menu is open, and the 'Assignments' option is highlighted with a red box. Below the menu, there is a banner for 'CSCI 699: Probabilistic and Generative Models'. On the left, there is a 'Slim Announcements Widget' containing two announcements: 'Lecture Slides Added and Note on Next Class' (Posted Wednesday, January 22, 2025 at 1:20 PM) and 'Welcome to Probabilistic and Generative Models!' (Posted Thursday, January 16, 2025 at 1:14 PM). There is also an 'Announcements' section with a single post from Willie Neiswanger. On the right, there are widgets for 'Calendar' (showing Friday, March 7, 2025 and Upcoming events), 'Activity Feed' (which is empty), and 'Latest Posts' (which is also empty). At the bottom, there is a 'Multi-Profile Widget'.

Today

Here's how to upload: (3) Click on any of the listed assignments.



The screenshot shows a course management system interface. At the top, the header includes the USC Viterbi School of Engineering logo, the course number 20251_30165 CSCI-699: Special To..., and various navigation links like Home, Announcements, Content, Activities, My Grades, Help, Course Tools, and Library Resources. On the right side of the header, there are icons for messaging, notifications, and user impersonation. Below the header, the main content area is titled "Assignments". A "View History" button is located above the assignment table. The assignment table has columns for Assignment, Completion Status, Score, and Evaluation Status. There are four rows in the table:

Assignment	Completion Status	Score	Evaluation Status
No Category			
Scribe Notes ⓘ	Not Submitted	- / 100	
Discussion Lead 1 - Questions ⓘ	Not Submitted	- / 100	
Discussion Lead 2 - Questions ⓘ	Not Submitted	- / 100	

A red rectangular box highlights the second row, which corresponds to the "Discussion Lead 1 - Questions" assignment. At the bottom right of the table, there is a "20 per page" dropdown menu.

Today

Here's how to upload: (4) Upload PDF (under *Add a File*)

The screenshot shows a Moodle assignment submission interface. At the top, the header includes the USC Viterbi School of Engineering logo, course information (20251_30165 CSCI-699: Special To...), and user details (Willie Neiswanger_Test Impersonating). The navigation bar below has links for Home, Announcements, Content, Activities, My Grades, Help, Course Tools, and Library Resources. The current page is 'Assignments > Discussion Lead 1 - Questions'. The main title is 'Discussion Lead 1 - Questions'. Below it, there is an 'Instructions' section with text about uploading questions for the 'Discussion Lead 1' role. A note says: 'After you have completed your discussion lead role (during the in-class paper presentation you have chosen), please upload a PDF file containing your ~5 discussion questions.' The 'Submit Assignment' section asks for files to submit, showing '(0) file(s) to submit'. It includes a note: 'After uploading, you must click Submit to complete the submission.' Three buttons are available: 'Add a File' (highlighted with a red box), 'Record Audio', and 'Record Video'. A 'Comments' section with rich text editing tools is also present. At the bottom, there are 'Submit' and 'Cancel' buttons.

Note – Next Week's Lecture

Note – Next Week's Lecture

Will be given by our very own Oliver Liu!



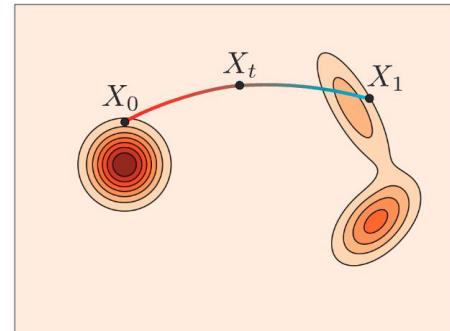
Oliver Liu
Teaching Assistant

Note – Next Week's Lecture

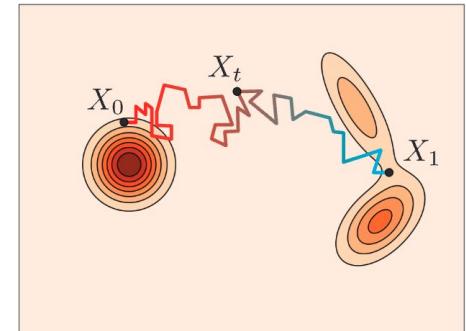


Will be given by our very own Oliver Liu!

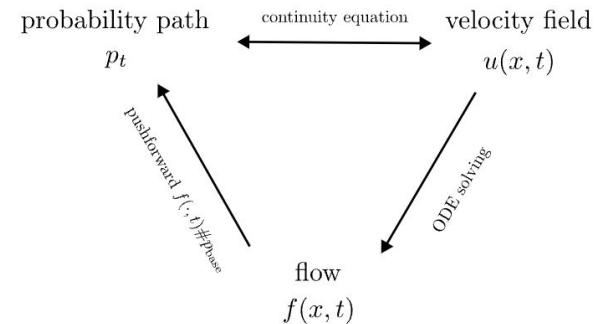
On the topic of **flow matching** ⇒



(a) Flow



(b) Diffusion



"Flow Matching Guide and Code",
Lipman et al., 2024
"A Visual Dive into Conditional Flow
Matching", 2024

Note – Next Week's Lecture

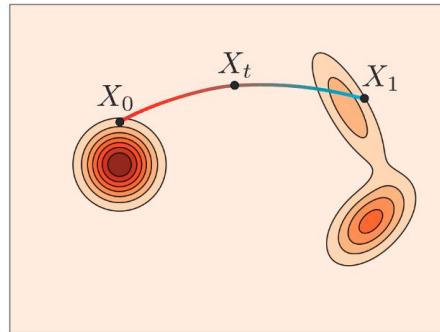


Will be given by our very own Oliver Liu!

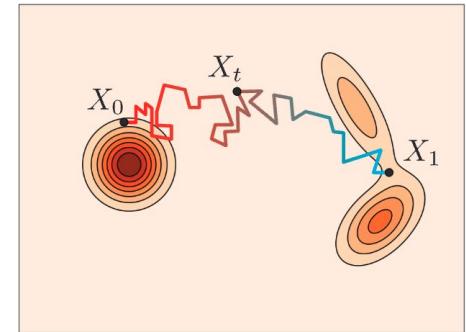
On the topic of **flow matching** ⇒

A good companion to the score-matching + diffusion lecture we are doing today.

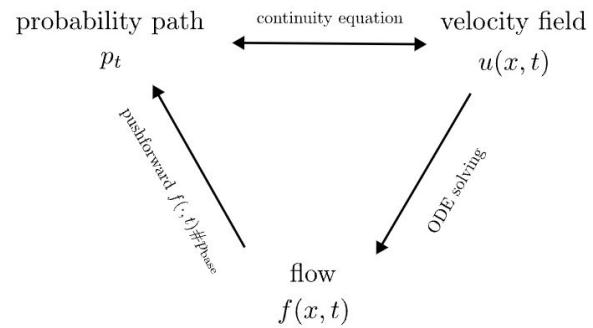
And will discuss another key deep generative modeling paradigm:
flow-based models.



(a) Flow



(b) Diffusion



"Flow Matching Guide and Code",
Lipman et al., 2024
"A Visual Dive into Conditional Flow Matching", 2024

Quick Review: VAE

Review – VAE Original 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

Review – VAE Original 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

Review – AEVB Algorithm

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!)

Review – 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 .

Review – 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:

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

Log of joint model
(unnormalized PDF)

Log of approximate
PDF

Review – The ELBO

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 .

Consider an ELBO for this model, written as:

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

Log of joint model
(unnormalized PDF)

Log of approximate
PDF

Note that the ELBO satisfies:

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

Evidence equals KL divergence plus ELBO.

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 $q_\phi(z)$.

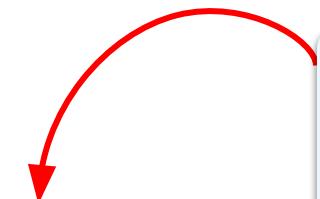
The Score Function Gradient Estimator

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**:

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



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 → 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

Expectation of the *log-likelihood*

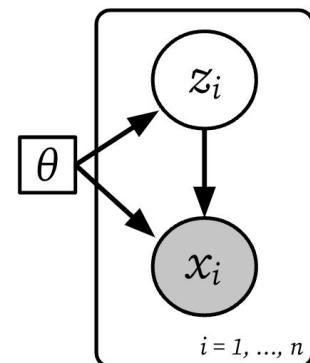
Reformulating the ELBO

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

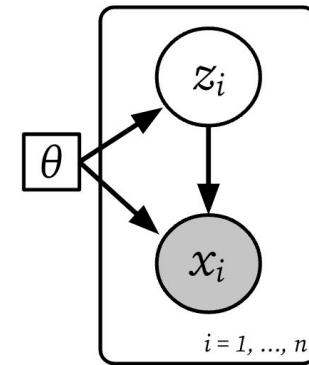
The reformulation of the ELBO is as follows:

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

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*

The first term:

- \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

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*

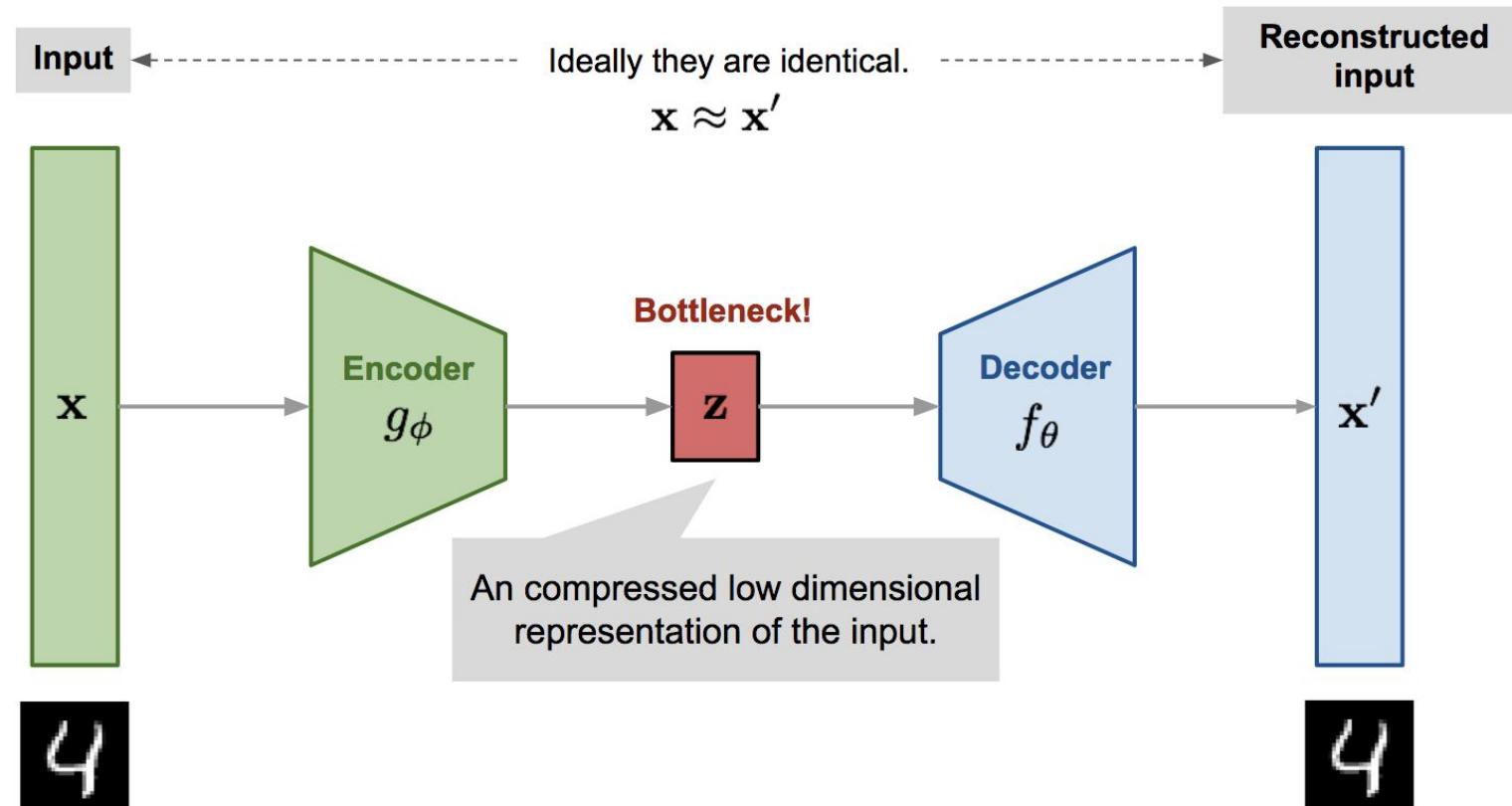
KL between *q posterior* and *p prior*

The first term:

- \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*.
- ... resembles a classic Autoencoder neural network.

Reconstruction Term \Leftrightarrow Autoencoder

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



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

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.

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) – 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 [11] 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 increasing noise levels in the score estimation process to close the loop. One of our first-order 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 plausible 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 shows generated samples from two different generative models. On the left, there are four images of diverse individuals from the CelebA-HQ dataset, each with a 4x4 grid of generated faces below it. On the right, there is a 4x4 grid of images from the CIFAR10 dataset, showing various objects like animals and vehicles.

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 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) 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 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., 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.

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

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), and maximum likelihood (Hyvonen, 2005), pass likelihood (Bengio, 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 (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 from closest point to the manifold. 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 $J(\theta, \phi)$ [15, 29, 11, 60].

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 images 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 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 (Giu et al., 2020), and

^{*}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 the 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

Peter 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 losy decomposition 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>.

arXiv:2006.11239v2 [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 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 losy decomposition 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

Diffenerative 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.

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
dirk@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 score-based generative model that can 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 models with millions of parameters 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 two errors 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-implicit ODEs (Song et al., 2019) on training data at each noise scale, and then uses Langevin dynamics to sample from the sequence of generated 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 can speed the DDPM training objective by explicitly computing 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; 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.

Denoising Diffusion Prob 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 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 annealing process to close the loop. The final model. One 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 [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 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 noisy decomposition 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

Diffenerative 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.

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 the score function with a sequence of unbiased gradient 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 distribution to make training tractable. For example, we expand the DDPM training objective to directly corrupt scores at coarse 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 (Cai 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 → [later/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

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 $\theta \in \mathbb{R}^d$.

⇒ We want to learn θ such that $p_\theta(x) \approx p(x)$.

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

To perform generative modeling, we want to learn model parameters θ .

Score-based Generative Models – Difficulties with MLE

To perform generative modeling, we want to learn model parameters θ .

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

To perform generative modeling, we want to learn model parameters θ .

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 $\tilde{p}_{\theta}(x)$.
- In VAEs, we perform MLE on θ via the use of an approximate posterior q .

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 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 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 data density $p(x)$, 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 data density $p(x)$, i.e.,

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

Gradient of the PDF with respect to x

Score-based Generative Models

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

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

Gradient of the PDF with respect to 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 data density $p(x)$, i.e.,

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

This score function can also be written:

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

Score-based Generative Models

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

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

This score function can also be written:

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

Definition of
unnormalized
PDF

Score-based Generative Models

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

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

This score function can also be written:

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

Second term is 0
(constant wrt x)

Definition of
unnormalized
PDF

Score-based Generative Models

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

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

This score function can also be written:

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

Second term is 0
(constant wrt x)

Definition of
unnormalized
PDF

Score-based Generative Models

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

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

This score function can also be written:

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

Also (in terms of energy based models):

Score-based Generative Models

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

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

This score function can also be written:

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

Also (in terms of energy based models):

Score-based Generative Models

First, we will aim to learn the **(Stein) score function** of the data density $p(x)$, 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

Also (in terms of energy
based models):

Score-based Generative Models

First, we will aim to learn the **(Stein) score function** of the data density $p(x)$, 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

Also (in terms of energy
based models):

Score-based Generative Models

Suppose we have observed a set of data: $D = \{x_1, x_2, \dots, x_n\} \stackrel{\text{iid}}{\sim} p(x)$.

And also suppose that we have (somehow) used D to learn an approximation $s_\theta(x)$ to the true score function:

Score-based Generative Models

Suppose we have observed a set of data: $D = \{x_1, x_2, \dots, x_n\} \stackrel{\text{iid}}{\sim} p(x)$.

And also suppose that we have (somehow) used D to learn 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 observed a set of data: $D = \{x_1, x_2, \dots, x_n\} \stackrel{\text{iid}}{\sim} p(x)$.

And also suppose that we have (somehow) used D to learn 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):

Score-based Generative Models

Suppose we have observed a set of data: $D = \{x_1, x_2, \dots, x_n\} \stackrel{\text{iid}}{\sim} p(x)$.

And also suppose that we have (somehow) used D to learn 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):

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

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

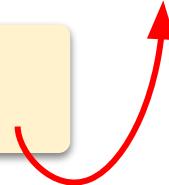
Score-based Generative Models

Therefore (after learning the score function):

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.



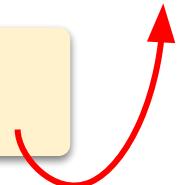
Score-based Generative Models

Therefore (after learning the score function):

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

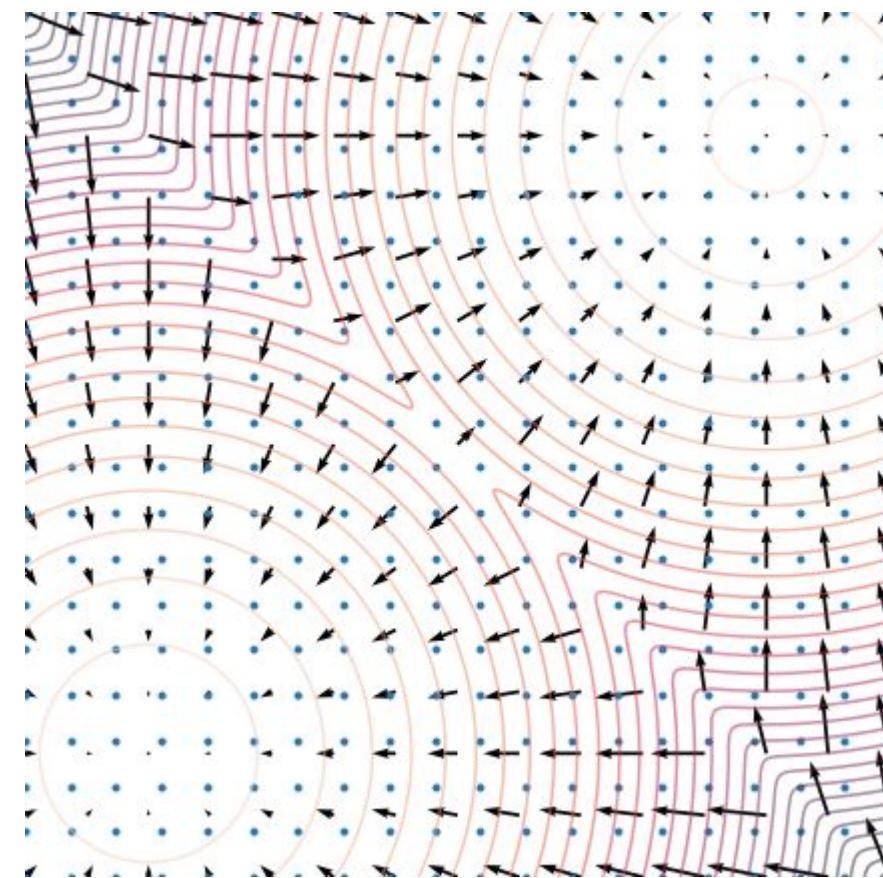
An illustration:

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.

(Many chains, randomly initialized).



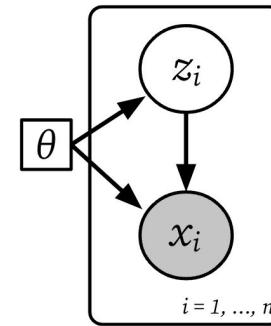
Score-based Generative Models – Intuition on Learning

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?

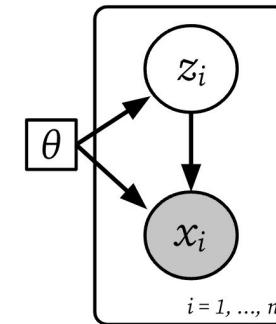
- In VAEs, we aim to learn a **direct map** from noise to data.
 - *I.e.*, directly learn the PDF of interest.
 - \Rightarrow a neural network that generates a sample from this PDF.



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 (along with noise) 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:

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

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:

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

Expectation taken over the data density.

Intuitively: this compares the squared L2 distance between the ground-truth data score function 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:

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

Expectation taken over the data density.

However, there is an **issue**: we don't know the ground truth PDF $p(x)$ nor 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.

Very clever trick!

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.

Score Matching

Get ready for some derivations :D

Start with the Fisher divergence:

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

Score Matching

Start with the Fisher divergence:

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

Definition of
score model

Score Matching

Start with the Fisher divergence:

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

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

Expand norm of the difference

$$= \mathbb{E}_{p(x)} [\|\nabla_x \log p(x)\|^2 - 2(\nabla_x \log p(x))^\top (\nabla_x \log p_\theta(x)) + \|\nabla_x \log p_\theta(x)\|^2]$$

Score Matching

Start with the Fisher divergence:

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

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

Expand norm of the difference

$$= \mathbb{E}_{p(x)} [\|\nabla_x \log p(x)\|^2 - 2(\nabla_x \log p(x))^\top (\nabla_x \log p_\theta(x)) + \|\nabla_x \log p_\theta(x)\|^2]$$

Can ignore when optimizing
for θ (is a constant)

Score Matching

Start with the Fisher divergence:

$$\begin{aligned}\mathcal{L}(\theta) &= \mathbb{E}_{p(x)} [\|\nabla_x \log p(x) - s_\theta(x)\|_2^2] \\ &= \mathbb{E}_{p(x)} [\|\nabla_x \log p(x) - \nabla_x \log p_\theta(x)\|_2^2] \\ &= \mathbb{E}_{p(x)} [\|\nabla_x \log p(x)\|^2 - 2(\nabla_x \log p(x))^\top (\nabla_x \log p_\theta(x)) + \|\nabla_x \log p_\theta(x)\|^2] \\ &= -2\mathbb{E}_{p(x)} [(\nabla_x \log p(x))^\top (\nabla_x \log p_\theta(x))] + \mathbb{E}_{p(x)} [\|\nabla_x \log p_\theta(x)\|^2] + C\end{aligned}$$

Rewriting without
first term

Score Matching

Start with the Fisher divergence:

$$\begin{aligned}\mathcal{L}(\theta) &= \mathbb{E}_{p(x)} [\|\nabla_x \log p(x) - s_\theta(x)\|_2^2] \\ &= \mathbb{E}_{p(x)} [\|\nabla_x \log p(x) - \nabla_x \log p_\theta(x)\|_2^2] \\ &= \mathbb{E}_{p(x)} [\|\nabla_x \log p(x)\|^2 - 2(\nabla_x \log p(x))^\top (\nabla_x \log p_\theta(x)) + \|\nabla_x \log p_\theta(x)\|^2] \\ &= -2\mathbb{E}_{p(x)} [(\nabla_x \log p(x))^\top (\nabla_x \log p_\theta(x))] + \underbrace{\mathbb{E}_{p(x)} [\|\nabla_x \log p_\theta(x)\|^2]}_{C} + C\end{aligned}$$

We can handle this term!
(via sampling from $p(x)$)

Score Matching

Start with the Fisher divergence:

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

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

$$= \mathbb{E}_{p(x)} [\|\nabla_x \log p(x)\|^2 - 2(\nabla_x \log p(x))^\top (\nabla_x \log p_\theta(x)) + \|\nabla_x \log p_\theta(x)\|^2]$$

$$= -2 \underbrace{\mathbb{E}_{p(x)} [(\nabla_x \log p(x))^\top (\nabla_x \log p_\theta(x))]}_{\text{Need to focus on how to handle this term!}} + \underbrace{\mathbb{E}_{p(x)} [\|\nabla_x \log p_\theta(x)\|^2]}_{\text{We can handle this term! (via sampling from } p(x) \text{)}} + C$$

Need to focus on how to handle this term!

We can handle this term!
(via sampling from $p(x)$)

Score Matching

Inside the expectation of this first term (consider the i -th component of vector):

$$\mathbb{E}_{p(x)} [\nabla_x \log p_\theta(x)_i \nabla_x \log p(x)_i]$$

Start with this
intractable term.

Score Matching

Inside the expectation of this first term (consider the i -th component of vector):

$$\begin{aligned} & \mathbb{E}_{p(x)} [\nabla_x \log p_\theta(x)_i \nabla_x \log p(x)_i] \\ &= \int_{-\infty}^{\infty} \nabla_x \log p_\theta(x)_i \nabla_x \log p(x)_i p(x) dx \end{aligned}$$

Write as an integral

Score Matching

Inside the expectation of this first term (consider the i -th component of vector):

$$\begin{aligned} & \mathbb{E}_{p(x)} [\nabla_x \log p_\theta(x)_i \nabla_x \log p(x)_i] \\ &= \int_{-\infty}^{\infty} \nabla_x \log p_\theta(x)_i \nabla_x \log p(x)_i p(x) dx \\ &= \int_{-\infty}^{\infty} \nabla_x \log p_\theta(x)_i \frac{\nabla_x p(x)_i}{p(x)} p(x) dx \end{aligned}$$

Apply chain rule for
gradient of a log

Score Matching

Inside the expectation of this first term (consider the i -th component of vector):

$$\begin{aligned} & \mathbb{E}_{p(x)} [\nabla_x \log p_\theta(x)_i \nabla_x \log p(x)_i] \\ &= \int_{-\infty}^{\infty} \nabla_x \log p_\theta(x)_i \nabla_x \log p(x)_i p(x) dx \\ &= \int_{-\infty}^{\infty} \nabla_x \log p_\theta(x)_i \frac{\nabla_x p(x)_i}{\cancel{p(x)}} \cancel{p(x)} dx \\ &= \int_{-\infty}^{\infty} \nabla_x \log p_\theta(x)_i \nabla_x p(x)_i dx \end{aligned}$$

Cancel out $p(x)$

Score Matching

Inside the expectation of this first term (consider the i -th component of vector):

$$\begin{aligned} & \mathbb{E}_{p(x)} [\nabla_x \log p_\theta(x)_i \nabla_x \log p(x)_i] \\ &= \int_{-\infty}^{\infty} \nabla_x \log p_\theta(x)_i \nabla_x \log p(x)_i p(x) dx \\ &= \int_{-\infty}^{\infty} \nabla_x \log p_\theta(x)_i \frac{\nabla_x p(x)_i}{\cancel{p(x)}} \cancel{p(x)} dx \\ &= \int_{-\infty}^{\infty} \nabla_x \log p_\theta(x)_i \nabla_x p(x)_i dx \end{aligned}$$

Doesn't look tractable or like an expectation... What to do here?

Score Matching

Inside the expectation of this first term (consider the i -th component of vector):

$$\begin{aligned} & \mathbb{E}_{p(x)} [\nabla_x \log p_\theta(x)_i \nabla_x \log p(x)_i] \\ &= \int_{-\infty}^{\infty} \nabla_x \log p_\theta(x)_i \nabla_x \log p(x)_i p(x) dx \\ &= \int_{-\infty}^{\infty} \nabla_x \log p_\theta(x)_i \frac{\nabla_x p(x)_i}{p(x)} \cancel{p(x)} dx \\ &= \int_{-\infty}^{\infty} \nabla_x \log p_\theta(x)_i \nabla_x p(x)_i dx \end{aligned}$$

Integration by parts ✨

Score Matching

Inside the expectation of this first term (consider the i -th component of vector):

$$\begin{aligned} & \mathbb{E}_{p(x)} [\nabla_x \log p_\theta(x)_i \nabla_x \log p(x)_i] \\ &= \int_{-\infty}^{\infty} \nabla_x \log p_\theta(x)_i \nabla_x \log p(x)_i p(x) dx \\ &= \int_{-\infty}^{\infty} \nabla_x \log p_\theta(x)_i \frac{\nabla_x p(x)_i}{p(x)} p(x) dx \quad \text{Integration by parts ✨} \\ &= \int_{-\infty}^{\infty} \nabla_x \log p_\theta(x)_i \nabla_x p(x)_i dx = - \int_{-\infty}^{\infty} (\nabla_x^2 \log p_\theta(x))_{i,i} p(x) dx \end{aligned}$$

(and some regularity assumptions)

Score Matching

Inside the expectation of this first term (consider the i -th component of vector):

$$\begin{aligned} & \mathbb{E}_{p(x)} [\nabla_x \log p_\theta(x)_i \nabla_x \log p(x)_i] \\ &= \int_{-\infty}^{\infty} \nabla_x \log p_\theta(x)_i \nabla_x \log p(x)_i p(x) dx \\ &= \int_{-\infty}^{\infty} \nabla_x \log p_\theta(x)_i \frac{\nabla_x p(x)_i}{p(x)} p(x) dx \quad \text{Integration by parts ✨} \\ &= \int_{-\infty}^{\infty} \nabla_x \log p_\theta(x)_i \nabla_x p(x)_i dx = - \int_{-\infty}^{\infty} (\nabla_x^2 \log p_\theta(x))_{i,i} p(x) dx \end{aligned}$$

i -th diagonal entry of
the Hessian matrix

(and some regularity assumptions)

Score Matching

Back to the Fisher divergence objective:

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

Score Matching

Back to the Fisher divergence objective:

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

From initial derivation

$$= -2\mathbb{E}_{p(x)} [(\nabla_x \log p(x))^\top (\nabla_x \log p_\theta(x))] + \mathbb{E}_{p(x)} [\|\nabla_x \log p_\theta(x)\|^2]$$

Score Matching

Back to the Fisher divergence objective:

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

From initial derivation

$$= -2\mathbb{E}_{p(x)} [(\nabla_x \log p(x))^\top (\nabla_x \log p_\theta(x))] + \mathbb{E}_{p(x)} [\|\nabla_x \log p_\theta(x)\|^2]$$

$$= 2 \mathbb{E}_{p(x)} [\underline{\text{tr}(\nabla_x^2 \log p_\theta(x))}] + \mathbb{E}_{p(x)} [\|\nabla_x \log p_\theta(x)\|_2^2]$$

Sub in first term

Score Matching

Back to the Fisher divergence objective:

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

From initial derivation

$$= -2\mathbb{E}_{p(x)} [(\nabla_x \log p(x))^\top (\nabla_x \log p_\theta(x))] + \mathbb{E}_{p(x)} [\|\nabla_x \log p_\theta(x)\|^2]$$

$$= 2\mathbb{E}_{p(x)} [\text{tr}(\nabla_x^2 \log p_\theta(x))] + \mathbb{E}_{p(x)} [\|\nabla_x \log p_\theta(x)\|_2^2]$$

Sub in first term

$$\propto \mathbb{E}_{p(x)} [\text{tr} (\nabla_x^2 \log p_\theta(x)) + \frac{1}{2} \|\nabla_x \log p_\theta(x)\|_2^2]$$

Can be rewritten like this

Score Matching

Fisher divergence – Monte Carlo estimate:

(Last line on previous slide)

$$\mathcal{L}(\theta) \propto \mathbb{E}_{p(x)} [\text{tr} (\nabla_x^2 \log p_\theta(x)) + \frac{1}{2} \|\nabla_x \log p_\theta(x)\|_2^2]$$

Score Matching

Fisher divergence – Monte Carlo estimate:

$$\mathcal{L}(\theta) \propto \mathbb{E}_{p(x)} [\text{tr} (\nabla_x^2 \log p_\theta(x)) + \frac{1}{2} \|\nabla_x \log p_\theta(x)\|_2^2]$$

⇒ Monte Carlo estimate!

Score Matching

Fisher divergence – Monte Carlo estimate:

$$\mathcal{L}(\theta) \propto \mathbb{E}_{p(x)} [\text{tr} (\nabla_x^2 \log p_\theta(x)) + \frac{1}{2} \|\nabla_x \log p_\theta(x)\|_2^2]$$

⇒ Monte Carlo estimate!

$$\approx \frac{1}{M} \sum_{m=1}^M \text{tr} (\nabla_x^2 \log p_\theta(x^{(m)})) + \frac{1}{2} \|\nabla_x \log p_\theta(x^{(m)})\|_2^2,$$

where $x^{(m)} \sim p(x)$

Score Matching

Fisher divergence – Monte Carlo estimate:

$$\mathcal{L}(\theta) \propto \mathbb{E}_{p(x)} [\text{tr} (\nabla_x^2 \log p_\theta(x)) + \frac{1}{2} \|\nabla_x \log p_\theta(x)\|_2^2]$$

⇒ Monte Carlo estimate!

$$\approx \frac{1}{M} \sum_{m=1}^M \text{tr} (\nabla_x^2 \log p_\theta(x^{(m)})) + \frac{1}{2} \|\nabla_x \log p_\theta(x^{(m)})\|_2^2,$$

where $x^{(m)} \sim p(x)$

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 numerical integration or sampling of the data to estimate the model parameters. 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 gradient even under mild conditions. The density function of the observed data does not appear in this formula, which simplifies the estimation. 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 $p(\cdot; \theta)$, where θ is an m -dimensional vector of parameters. We want to estimate the parameter θ from \mathbf{x} , i.e. we want to approximate $p_\mathbf{x}(\cdot)$ by $p(\cdot; \hat{\theta})$ for the estimated parameter value $\hat{\theta}$. (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 $Z(\theta)$:

$$p(\xi; \theta) = \frac{1}{Z(\theta)} q(\xi; \theta).$$

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.

Score Matching

Fisher divergence – Monte Carlo estimate:

$$\mathcal{L}(\theta) \propto \mathbb{E}_{p(x)} [\text{tr} (\nabla_x^2 \log p_\theta(x)) + \frac{1}{2} \|\nabla_x \log p_\theta(x)\|_2^2]$$

⇒ Monte Carlo estimate!

$$\approx \frac{1}{M} \sum_{m=1}^M \text{tr} (\nabla_x^2 \log p_\theta(x^{(m)})) + \frac{1}{2} \|\nabla_x \log p_\theta(x^{(m)})\|_2^2,$$

where $x^{(m)} \sim p(x)$

Important: scalability of method, and how to actually optimize this loss to learn score network $s_\theta(x)$:

- Hard to scale up equation below to high dimensions.
- More scalable methods: *sliced score matching*, and *denoising score matching*.

Summary

To summarize the (naive/vanilla) procedure:

Summary

To summarize the (naive/vanilla) procedure:

- (1) Use Monte Carlo estimate from previous slide to learn a score network $s_\theta(x)$.

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

Summary

To summarize the (naive/vanilla) procedure:

- (1) Use Monte Carlo estimate from previous slide to learn a score network $s_\theta(x)$.

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

- (2) Generate samples via Langevin Monte Carlo

Summary

To summarize the (naive/vanilla) procedure:

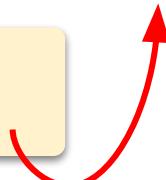
- (1) Use Monte Carlo estimate from previous slide to learn a score network $s_\theta(x)$.

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

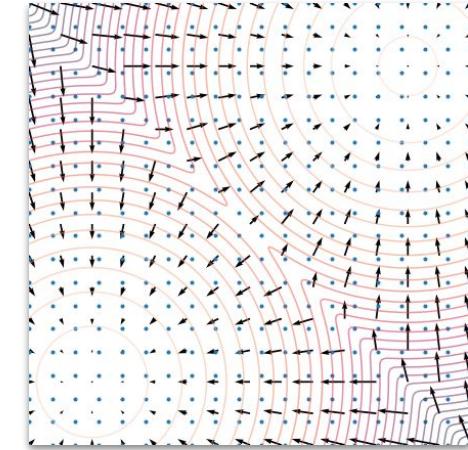
- (2) Generate samples via Langevin Monte Carlo

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

Note the learned score function is swapped in.



Summary



To summarize the (naive/vanilla) procedure:

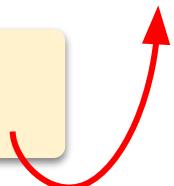
- (1) Use Monte Carlo estimate from previous slide to learn a score network $s_\theta(x)$.

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

- (2) Generate samples via Langevin Monte Carlo

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

Note the learned score
function is swapped in.



Difficulties to Solve in Practice

Unfortunately, this naive/vanilla procedure has some pitfalls in practice!

⇒ Running this without modification does not work great.

Difficulties to Solve in Practice

Unfortunately, this naive/vanilla procedure has some pitfalls in practice!

⇒ Running this without modification does not work great.

Main issue: estimate score function is inaccurate in low-density regions of $p(x)$.

Difficulties to Solve in Practice

Unfortunately, this naive/vanilla procedure has some pitfalls in practice!

⇒ Running this without modification does not work great.

Main issue: estimate score function is inaccurate in low-density regions of $p(x)$.

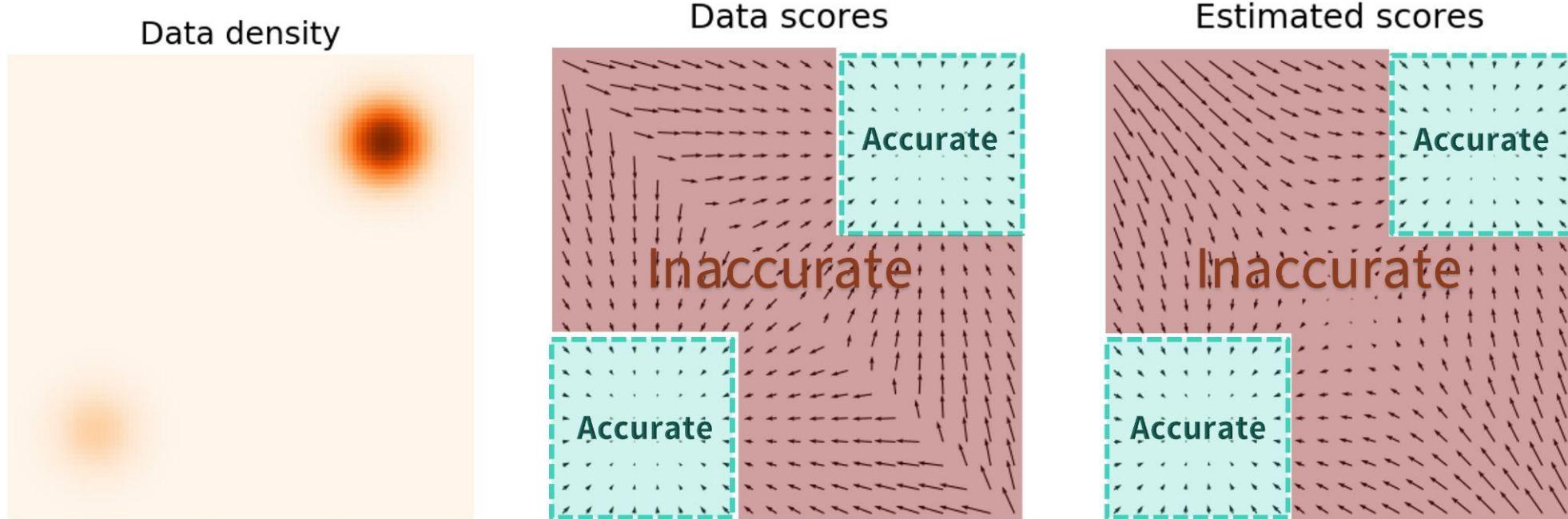
Intuition – to learn the score function, we minimized the Fisher divergence:

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

Note the loss is only large
when $p(x)$ has a high value!

⇒ When taking Monte Carlo samples from
 $p(x)$, we only get samples in high-mass areas

Difficulties to Solve in Practice



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

Note the loss is only large
when $p(x)$ has a high value!

⇒ When taking Monte Carlo samples from $p(x)$, we only get samples in high-mass areas

Source: Yang

Difficulties to Solve in Practice

Consequences of this:

Difficulties to Solve in Practice

Consequences of this:

- Our initial samples will likely be in low-density regions (especially when our data is high dimensional, e.g., images.)

Difficulties to Solve in Practice

Consequences of this:

- Our initial samples will likely be in low-density regions (especially when our data is high dimensional, e.g., images.)
- ⇒ Langevin Monte Carlo becomes “derailed” at the beginning of the procedure.

Difficulties to Solve in Practice

Consequences of this:

- Our initial samples will likely be in low-density regions (especially when our data is high dimensional, e.g., images.)
- ⇒ Langevin Monte Carlo becomes “derailed” at the beginning of the procedure.
- ⇒ Unable to generate high quality samples.

Difficulties to Solve in Practice

However...

- Suppose that we perturb the data points with noise, and train score-based models on the noisy data points instead.

Difficulties to Solve in Practice

However...

- Suppose that we perturb the data points with noise, and train score-based models on the noisy data points instead.
- When the noise is large enough, it can populate low-density regions of $p(x)$.

Difficulties to Solve in Practice

However...

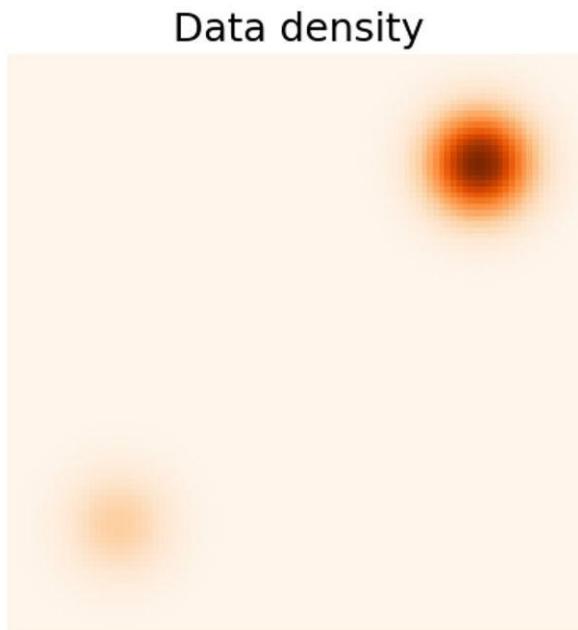
- Suppose that we perturb the data points with noise, and train score-based models on the noisy data points instead.
- When the noise is large enough, it can populate low-density regions of $p(x)$.
- ⇒ This improves the accuracy of estimated scores (for a slightly different/biased distribution).

Difficulties to Solve in Practice

Visually, using previous figure:

Difficulties to Solve in Practice

Visually, using previous figure:



Before: No perturbation.

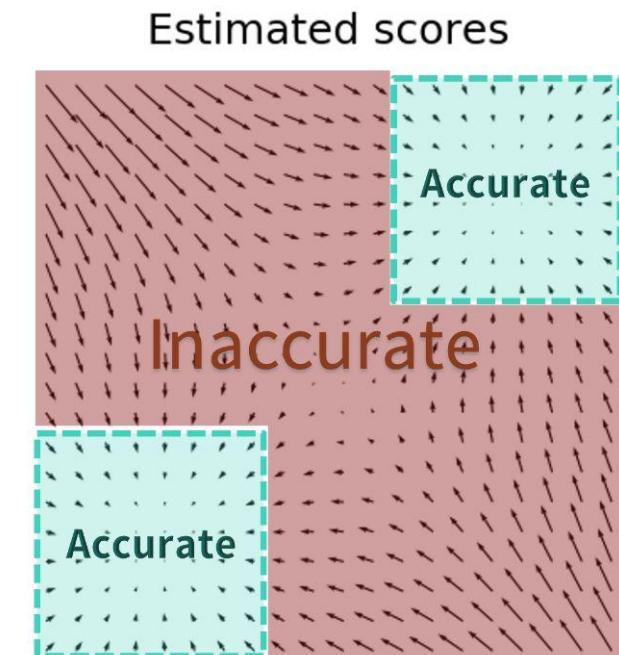
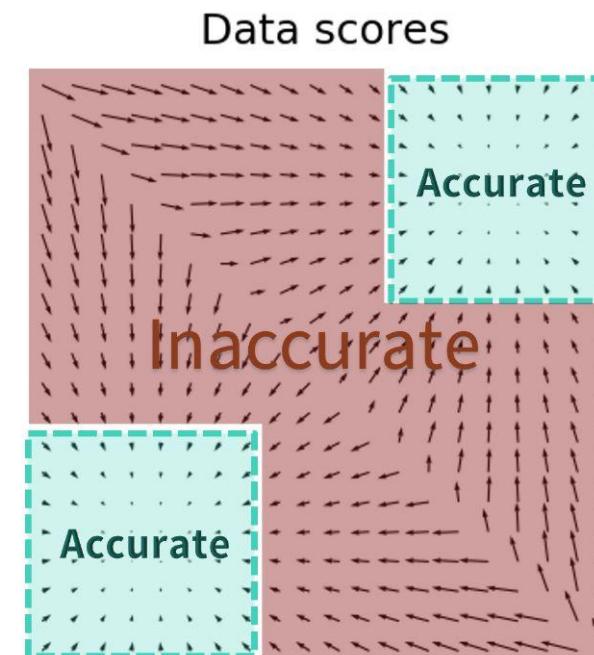
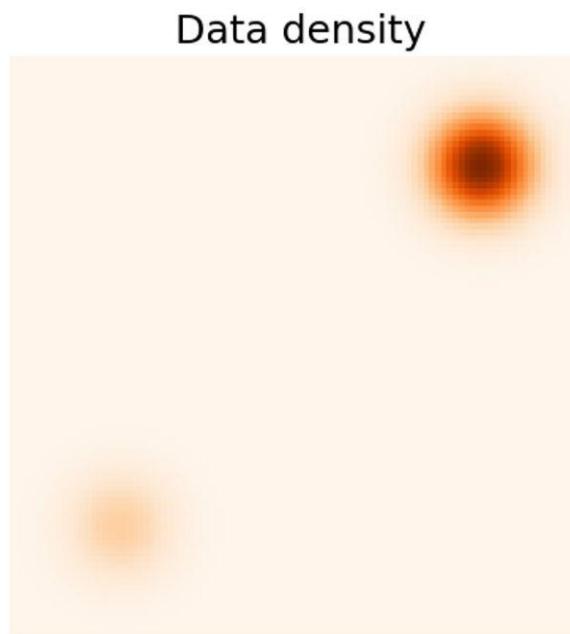
⇒ Estimated score is inaccurate in low-density regions.

Difficulties to Solve in Practice

Visually, using previous figure:

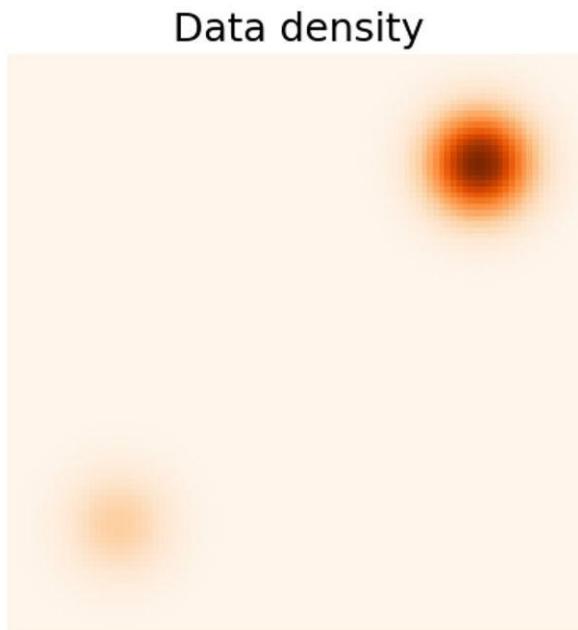
Before: No perturbation.

⇒ Estimated score is inaccurate in low-density regions.



Difficulties to Solve in Practice

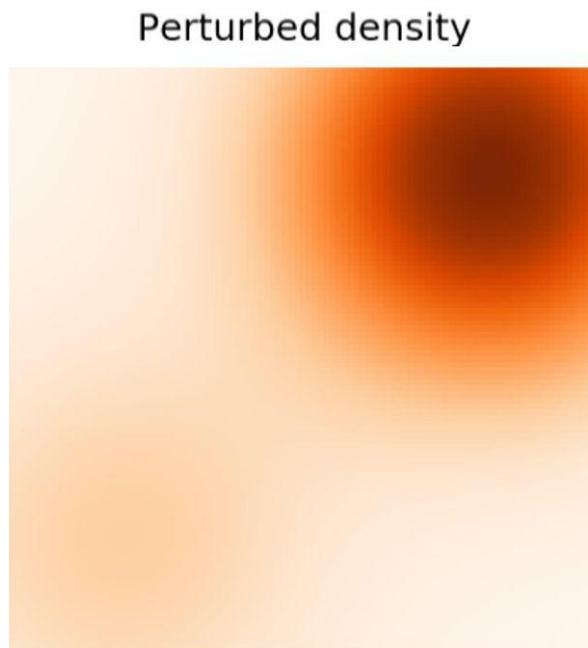
Visually, using previous figure:



Instead: **Perturb (add noise to) $p(x)$.**
⇒ Produces more samples in low-density regions, thus making score more accurate.

Difficulties to Solve in Practice

Visually, using previous figure:

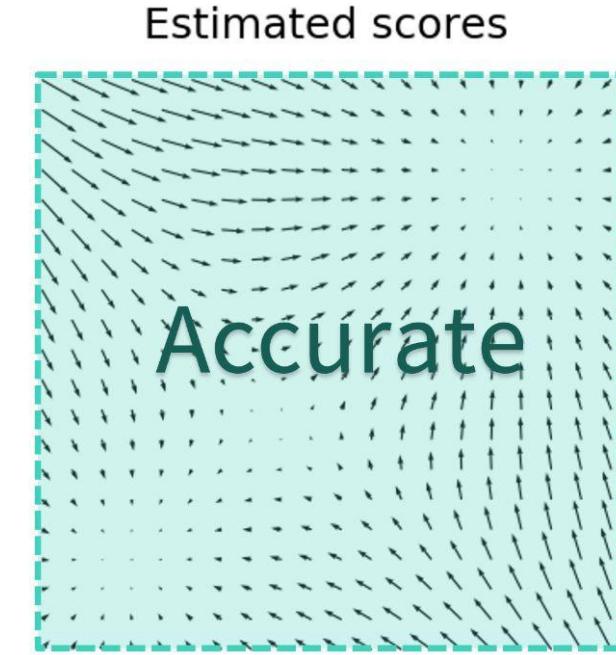
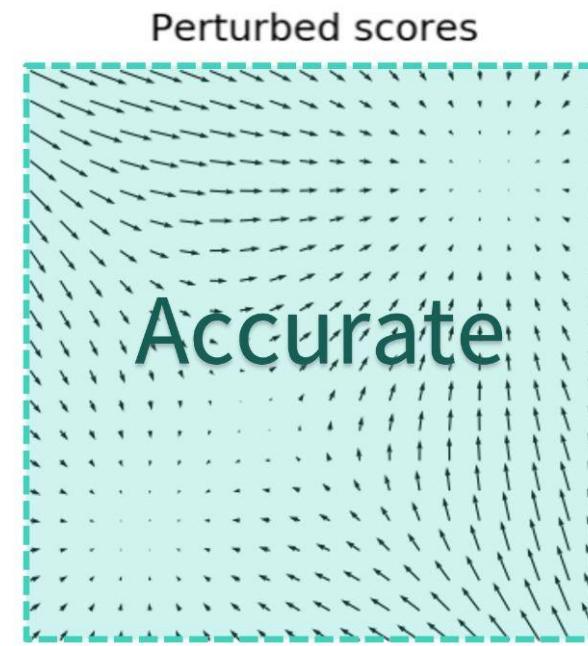
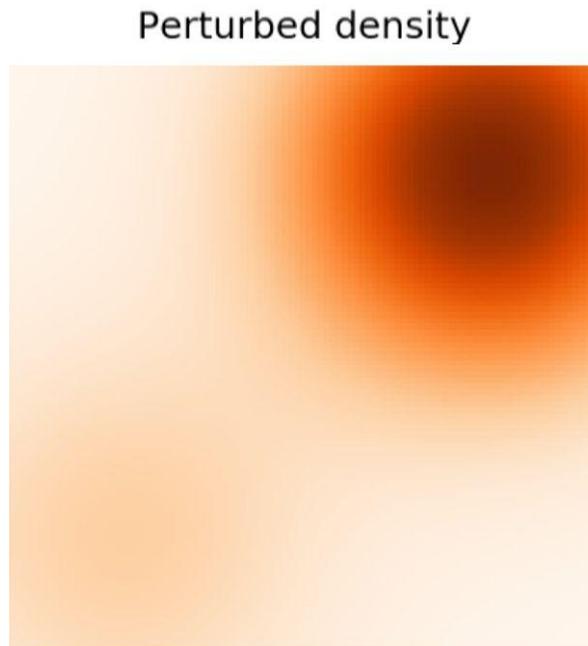


Instead: **Perturb (add noise to) $p(x)$.**
⇒ Produces more samples in low-density regions, thus making score more accurate.

Difficulties to Solve in Practice

Visually, using previous figure:

Instead: Perturb (add noise to) $p(x)$.
⇒ Produces more samples in low-density regions, thus making score more accurate.



Score-based Generative Models – Multiple Noise Perturbations

So we have a tradeoff:

Score-based Generative Models – Multiple Noise Perturbations

So we have a tradeoff:

Decreasing the noise makes score matching **harder** (especially at start of LMC), but yields more-accurate (*unbiased*) scores.

Score-based Generative Models – Multiple Noise Perturbations

So we have a tradeoff:

Decreasing the noise makes score matching **harder** (especially at start of LMC), but yields more-accurate (*unbiased*) scores.

Increasing the noise makes score matching **easier** (samples cover low-density regions), but yields less-accurate (*biased*) scores.

Score-based Generative Models – Multiple Noise Perturbations

So we have a tradeoff:

Decreasing the noise makes score matching **harder** (especially at start of LMC), but yields more-accurate (*unbiased*) scores.

Increasing the noise makes score matching **easier** (samples cover low-density regions), but yields less-accurate (*biased*) scores.

Solution:

Perform score matching (learn the score function) for *multiple noise levels simultaneously!*

Score-based Generative Models – Multiple Noise Perturbations

So we have a tradeoff:

Decreasing the noise makes score matching **harder** (especially at start of LMC), but yields more-accurate (*unbiased*) scores.

Increasing the noise makes score matching **easier** (samples cover low-density regions), but yields less-accurate (*biased*) scores.

Solution:

Perform score matching (learn the score function) for *multiple noise levels simultaneously!*

And then, to sample in LMC, we will traverse through the noise levels (from high to low noise).

Score-based Generative Models – Multiple Noise Perturbations

In particular:

Score-based Generative Models – Multiple Noise Perturbations

In particular:

- Suppose we perturb the data with isotropic Gaussian noise.

Score-based Generative Models – Multiple Noise Perturbations

In particular:

- Suppose we perturb the data with isotropic Gaussian noise.
- And that we have L increasing noise levels: $i = 1, \dots, L$.

Score-based Generative Models – Multiple Noise Perturbations

In particular:

- Suppose we perturb the data with isotropic Gaussian noise.
- And that we have L increasing noise levels: $i = 1, \dots, L$.
- For each noise level i , we can add Gaussian noise $\mathcal{N}(0, \sigma_i^2 I)$ to $p(x)$ to form the *noise-perturbed distribution*:

Score-based Generative Models – Multiple Noise Perturbations

In particular:

- Suppose we perturb the data with isotropic Gaussian noise.
- And that we have L increasing noise levels: $i = 1, \dots, L$.
- For each noise level i , we can add Gaussian noise $\mathcal{N}(0, \sigma_i^2 I)$ to $p(x)$ to form the *noise-perturbed distribution*:

$$p_{\sigma_i}(x) = \int p(x') \mathcal{N}(x \mid x', \sigma_i^2 I) dx' = \mathbb{E}_{p(x')} [\mathcal{N}(x \mid x', \sigma_i^2 I)]$$

Score-based Generative Models – Multiple Noise Perturbations

Next, we can estimate the score function for each noise-perturbed distribution i , which we can denote $s_\theta(x, i)$, using score matching.

Score-based Generative Models – Multiple Noise Perturbations

Next, we can estimate the score function for each noise-perturbed distribution i , which we can denote $s_\theta(x, i)$, using score matching.

At the end of this procedure, our score network should approximate:

$$s_\theta(x, i) \approx \nabla_x \log p_{\sigma_i}(x), \quad i = 1, 2, \dots, L$$

Score-based Generative Models – Multiple Noise Perturbations

Next, we can estimate the score function for each noise-perturbed distribution i , which we can denote $s_\theta(x, i)$, using score matching.

At the end of this procedure, our score network should approximate:

$$s_\theta(x, i) \approx \nabla_x \log p_{\sigma_i}(x), \quad i = 1, 2, \dots, L$$

We will call this model the **noise conditional score network (NCSN)**.

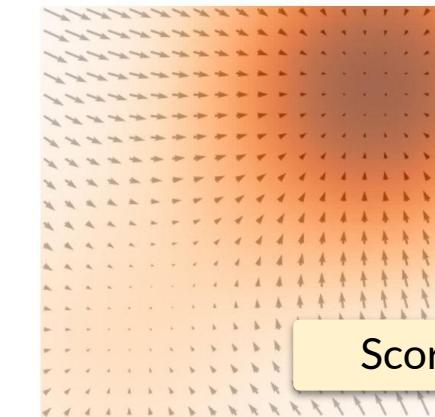
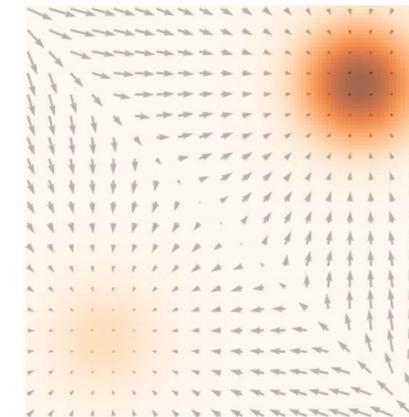
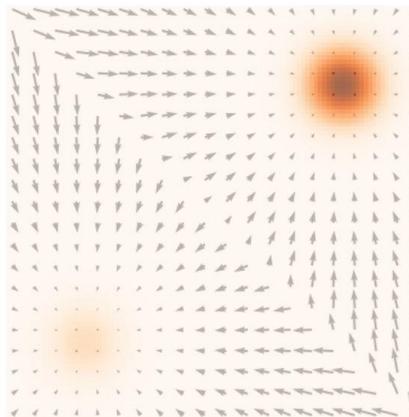
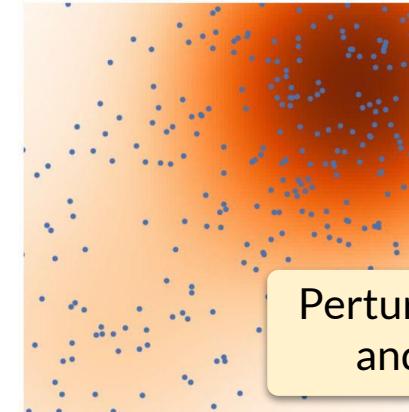
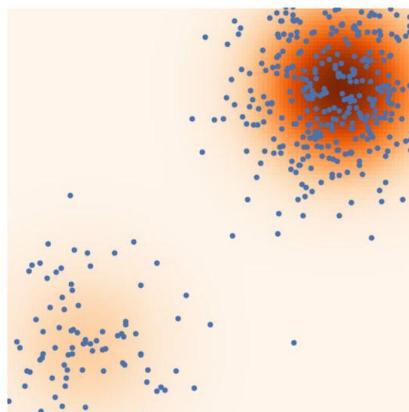
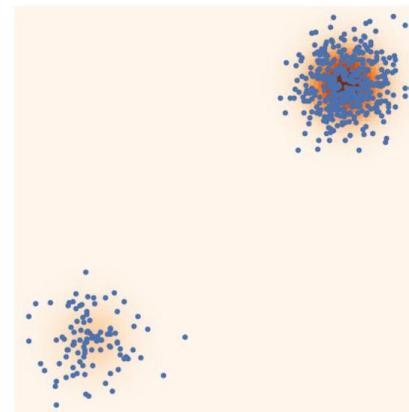
Score-based Generative Models – Multiple Noise Perturbations

Visualizing
this:

Score-based Generative Models – Multiple Noise Perturbations

Visualizing
this:

$$\sigma_1 < \sigma_2 < \sigma_3$$



Increasing noise level →

Perturbed density,
and samples

Score function

Score-based Generative Models – Multiple Noise Perturbations

If we perturb a data with multiple scales of Gaussian noise, we see an increasingly noisy data point.

An image example:



Note: it is easy for us to sample this perturbed data (given our original dataset).

Score-based Generative Models – Learning the NCSN

To learn the noise-conditional score network (NCSN) $s_\theta(x, i)$:

Score-based Generative Models – Learning the NCSN

To learn the noise-conditional score network (NCSN) $s_\theta(x, i)$:

We use a modified training objective – a weighted sum of Fisher divergences over all of the noise scales:

$$\mathcal{L}(\theta) = \sum_{i=1}^L \lambda(i) \mathbb{E}_{p_{\sigma_i}(x)} [\|\nabla_x \log p_{\sigma_i}(x) - s_\theta(x, i)\|_2^2]$$

Score-based Generative Models – Learning the NCSN

To learn the noise-conditional score network (NCSN) $s_\theta(x, i)$:

We use a modified training objective – a weighted sum of Fisher divergences over all of the noise scales:

$$\mathcal{L}(\theta) = \sum_{i=1}^L \lambda(i) \mathbb{E}_{p_{\sigma_i}(x)} [\|\nabla_x \log p_{\sigma_i}(x) - s_\theta(x, i)\|_2^2]$$

Where the weighting function is often taken to be $\lambda(i) = \sigma_i^2$.

Score-based Generative Models – Generation from NCSN

Finally, to sample from this model:

Score-based Generative Models – Generation from NCSN

Finally, to sample from this model:

Similar to before, run Langevin Monte Carlo (LMC) using the learned (noise-conditional) score network $s_\theta(x, i)$...

While sequentially stepping (backwards) through the noise levels: $i = L, L-1, L-2, \dots, 1$

Score-based Generative Models – Generation from NCSN

Finally, to sample from this model:

Similar to before, run Langevin Monte Carlo (LMC) using the learned (noise-conditional) score network $s_\theta(x, i)$...

While sequentially stepping (backwards) through the noise levels: $i = L, L-1, L-2, \dots, 1$

⇒ This is LMC with an annealed noise level/step size! (Mentioned in MCMC lecture).

Score-based Generative Models – Generation from NCSN

Altogether, how does this look?

Score-based Generative Models – Generation from NCSN

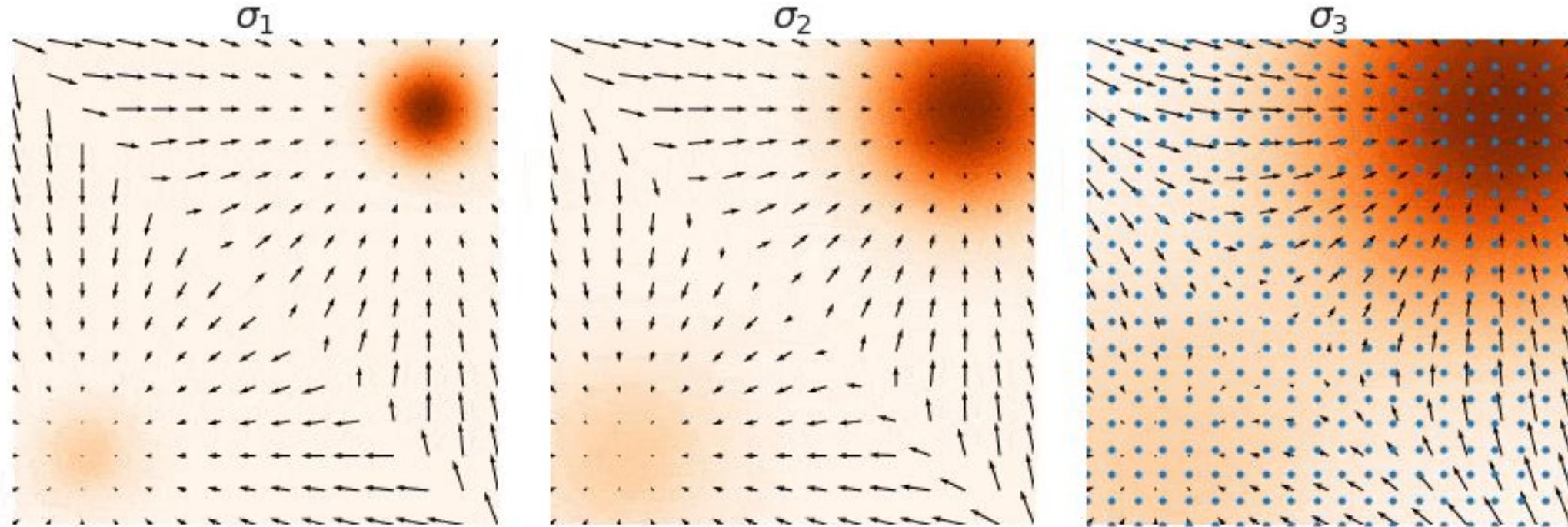
Altogether, how does this look?



Score-based Generative Models – Generation from NCSN

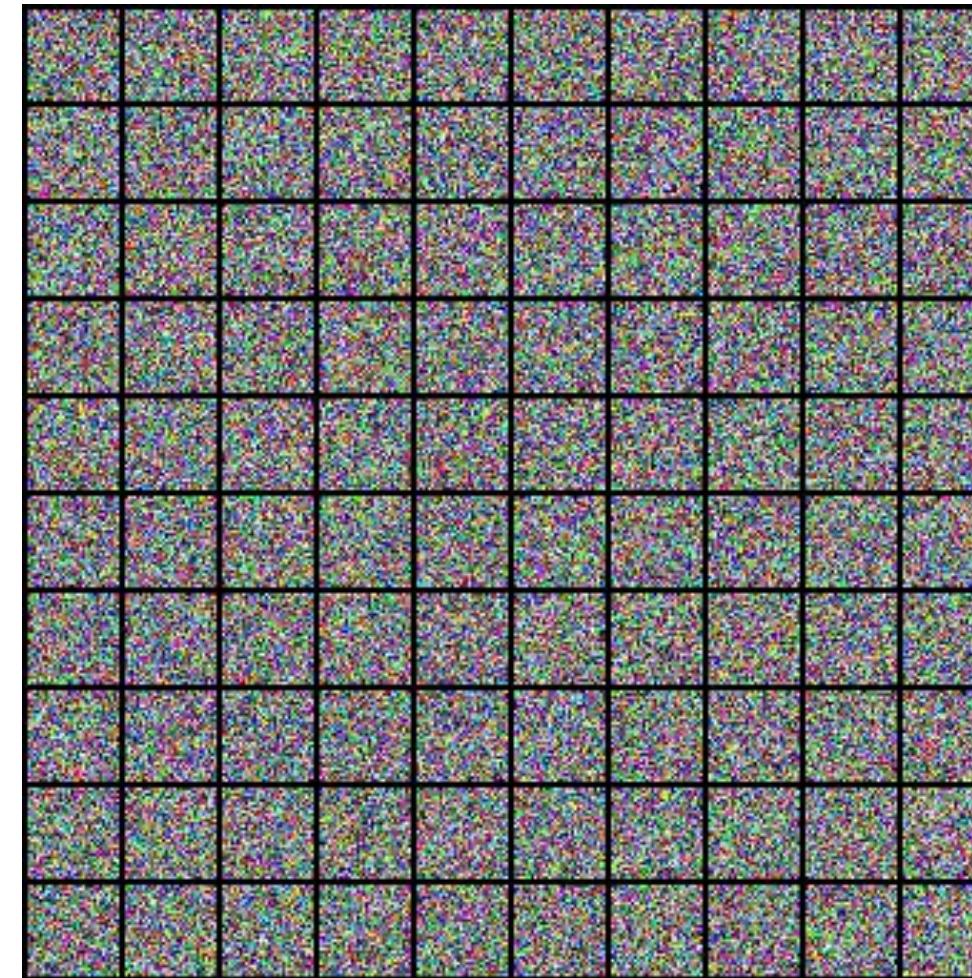
Altogether, how does this look?

← Proceeding from **right to left**.



Score-based Generative Models – Generation from NCSN

And on higher dimensional data
(face images):



Connections to Diffusion Models

Connections to Diffusion Models

You can see how score-based models start to look like what are typically called **diffusion models**.

Connections to Diffusion Models

You can see how score-based models start to look like what are typically called **diffusion models**.

Diffusion models can be viewed in the following way:

Connections to Diffusion Models

You can see how score-based models start to look like what are typically called **diffusion models**.

Diffusion models can be viewed in the following way:

- Iteratively add noise to data points (via a *forward diffusion process*).

Connections to Diffusion Models

You can see how score-based models start to look like what are typically called **diffusion models**.

Diffusion models can be viewed in the following way:

- Iteratively add noise to data points (via a *forward diffusion process*).
- Then learn a *reverse diffusion process*, which iteratively denoises data.

Connections to Diffusion Models

You can see how score-based models start to look like what are typically called **diffusion models**.

Diffusion models can be viewed in the following way:

- Iteratively add noise to data points (via a *forward diffusion process*).
- Then learn a *reverse diffusion process*, which iteratively denoises data.
- By sampling from the reverse diffusion process, can generate new data samples.

Connections to Diffusion Models

A popular example is the
denoising diffusion
probabilistic model (DDPM),
NeurIPS 2019.

(We will hear a presentation of
this paper later today.)

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 best results are obtained by training on a weighted variational bound designed according to a novel connection between diffusion 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 autoregressive decoding. On the unconditional CIFAR10 dataset, we obtain an Inception score of 9.46 and a state-of-the-art FID score of 3.17. On 256x256 LSUN, we obtain sample quality similar to ProgressiveGAN. Our implementation is available at <https://github.com/jonathanho/diffusion>.

1 Introduction

Deep 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 autoencoders (VAEs) have synthesized striking image and audio samples [14, 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)

Connections to Diffusion Models

But for a brief overview:

Connections to Diffusion Models

But for a brief overview:

- Define the DDPM as a probabilistic model.
- Where we view the *forward diffusion process* as a variational approximation to the posterior.

Connections to Diffusion Models

But for a brief overview:

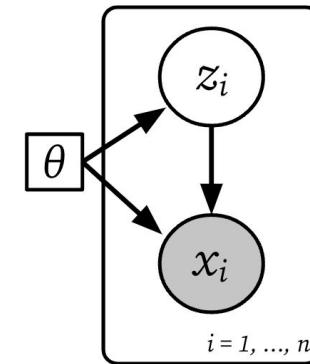
- Define the DDPM as a probabilistic model.
- Where we view the *forward diffusion process* as a variational approximation to the posterior.
- And we view the *reverse diffusion process* as a Markov chain, with generative model parameters θ that we'd like to learn.

Connections to Diffusion Models

But for a brief overview:

- Define the DDPM as a probabilistic model.
- Where we view the *forward diffusion process* as a variational approximation to the posterior.
- And we view the *reverse diffusion process* as a Markov chain, with generative model parameters θ that we'd like to learn.

Somewhat similar
in this way to VAEs

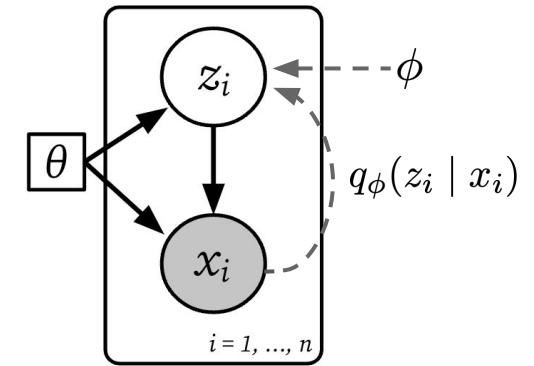


Connections to Diffusion Models

But for a brief overview:

- Define the DDPM as a probabilistic model.
- Where we view the *forward diffusion process* as a variational approximation to the posterior.
- And we view the *reverse diffusion process* as a Markov chain, with generative model parameters θ that we'd like to learn.

Somewhat similar
in this way to VAEs

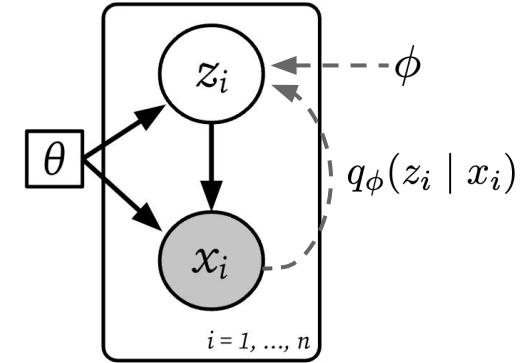


Connections to Diffusion Models

But for a brief overview:

- Define the DDPM as a probabilistic model.
- Where we view the *forward diffusion process* as a variational approximation to the posterior.
- And we view the *reverse diffusion process* as a Markov chain, with generative model parameters θ that we'd like to learn.
- And similar to VAE models, we maximize the ELBO while learning θ .

Somewhat similar
in this way to VAEs



Connections to Diffusion Models

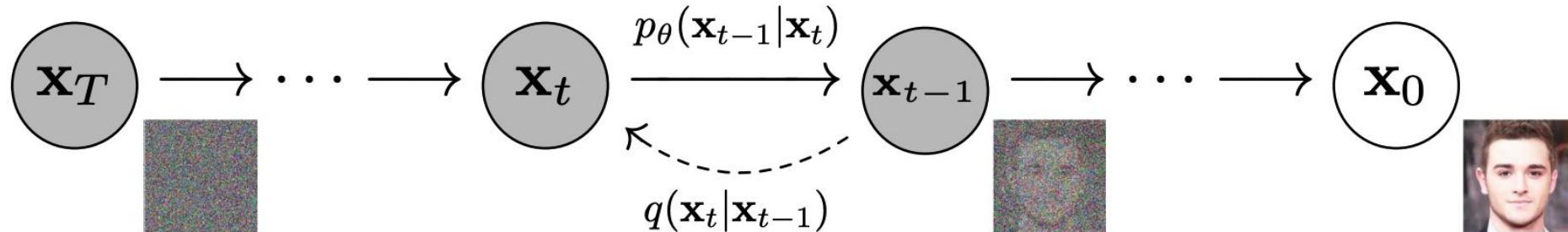
But for a brief overview:

The graphical model for the DDPM is popularly illustrated as:

Connections to Diffusion Models

But for a brief overview:

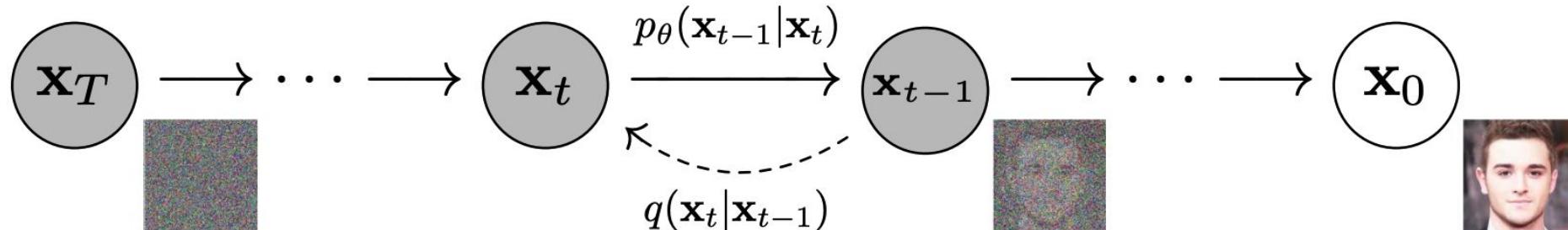
The graphical model for the DDPM is popularly illustrated as:



Connections to Diffusion Models

But for a brief overview:

The graphical model for the DDPM is popularly illustrated as:



Notice q is not learned (like in VAEs), but is a fixed process
... except for noise level parameters, in some cases.

Connections to Diffusion Models

So, to be explicit, the probabilistic model (Bayesian network) is:

$$p_{\theta}(x_0, \dots, x_T) = p_{\theta}(x_T) \prod_{t=1}^T p_{\theta}(x_{t-1} \mid x_t)$$

Connections to Diffusion Models

So, to be explicit, the probabilistic model (Bayesian network) is:

$$p_{\theta}(x_0, \dots, x_T) = p_{\theta}(x_T) \prod_{t=1}^T p_{\theta}(x_{t-1} \mid x_t)$$

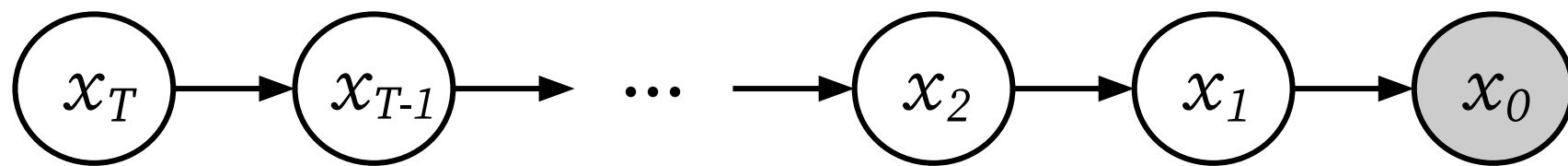
A Markov chain (with one node observed)

Connections to Diffusion Models

So, to be explicit, the probabilistic model (Bayesian network) is:

$$p_{\theta}(x_0, \dots, x_T) = p_{\theta}(x_T) \prod_{t=1}^T p_{\theta}(x_{t-1} \mid x_t)$$

A Markov chain (with one node observed)



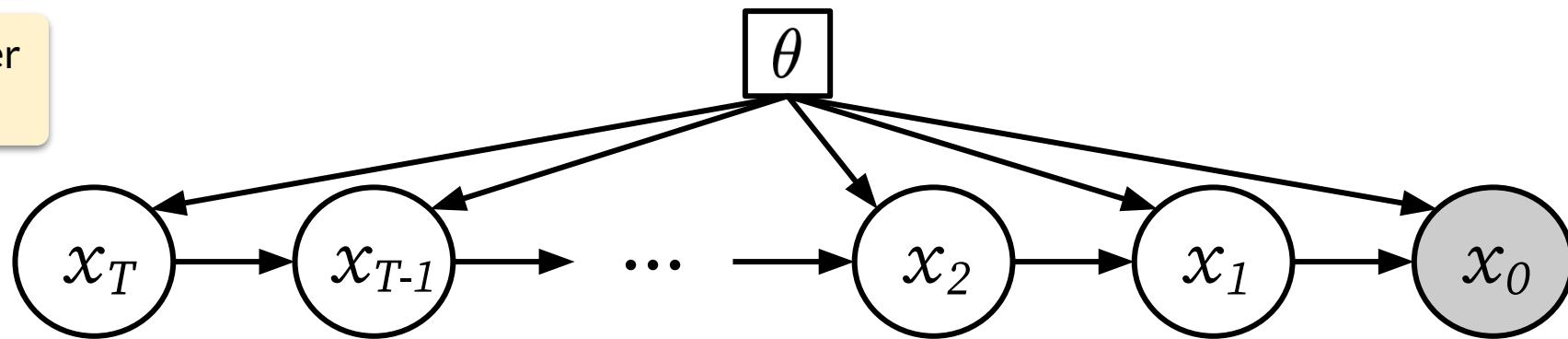
Connections to Diffusion Models

So, to be explicit, the probabilistic model (Bayesian network) is:

$$p_{\theta}(x_0, \dots, x_T) = p_{\theta}(x_T) \prod_{t=1}^T p_{\theta}(x_{t-1} \mid x_t)$$

A Markov chain (with one node observed)

With parameter θ shown



Connections to Diffusion Models

While the variational family of *posterior approximations* can be written:

Connections to Diffusion Models

While the variational family of *posterior approximations* can be written:

$$q(x_1, \dots, x_T \mid x_0) = \prod_{t=1}^T q(x_t \mid x_{t-1})$$

Connections to Diffusion Models

While the variational family of *posterior approximations* can be written:

$$q(x_1, \dots, x_T \mid x_0) = \prod_{t=1}^T q(x_t \mid x_{t-1})$$

(Another Markov chain!)

Connections to Diffusion Models

While the variational family of *posterior approximations* can be written:

$$q(x_1, \dots, x_T \mid x_0) = \prod_{t=1}^T q(x_t \mid x_{t-1})$$

(Another Markov chain!)

The graphical model for this distribution might be drawn as:

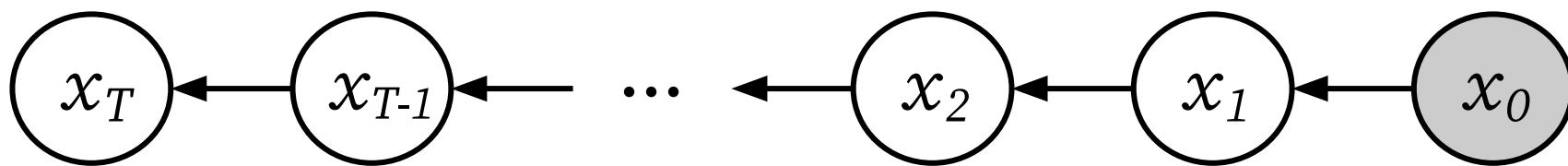
Connections to Diffusion Models

While the variational family of *posterior approximations* can be written:

$$q(x_1, \dots, x_T \mid x_0) = \prod_{t=1}^T q(x_t \mid x_{t-1})$$

(Another Markov chain!)

The graphical model for this distribution might be drawn as:



Connections to Diffusion Models

While the variational family of *posterior approximations* can be written:

$$q(x_1, \dots, x_T \mid x_0) = \prod_{t=1}^T q(x_t \mid x_{t-1})$$

(Another Markov chain!)

And sticking both together we get:

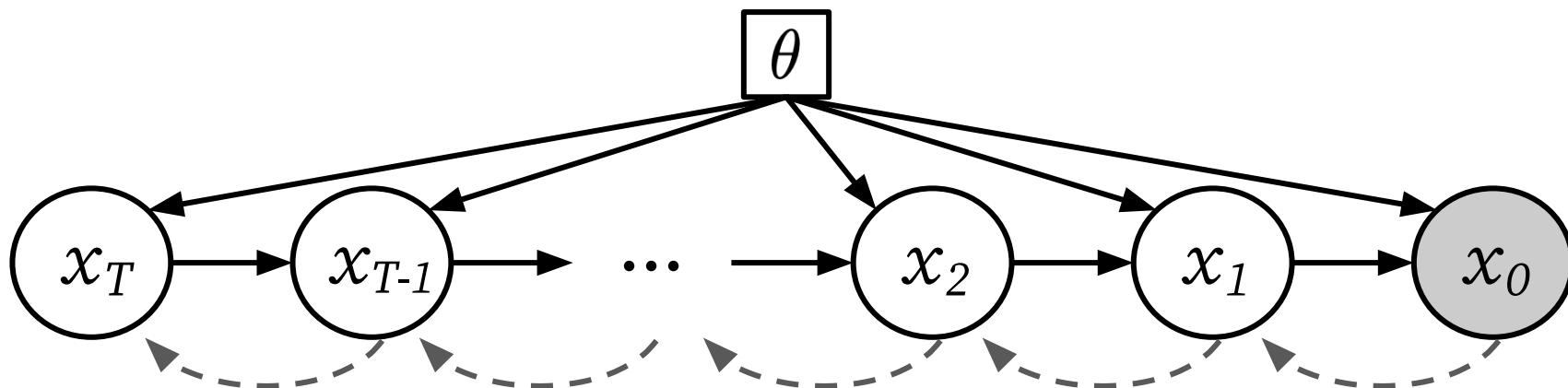
Connections to Diffusion Models

While the variational family of *posterior approximations* can be written:

$$q(x_1, \dots, x_T \mid x_0) = \prod_{t=1}^T q(x_t \mid x_{t-1})$$

(Another Markov chain!)

And sticking both together we get:



Connections to Diffusion Models

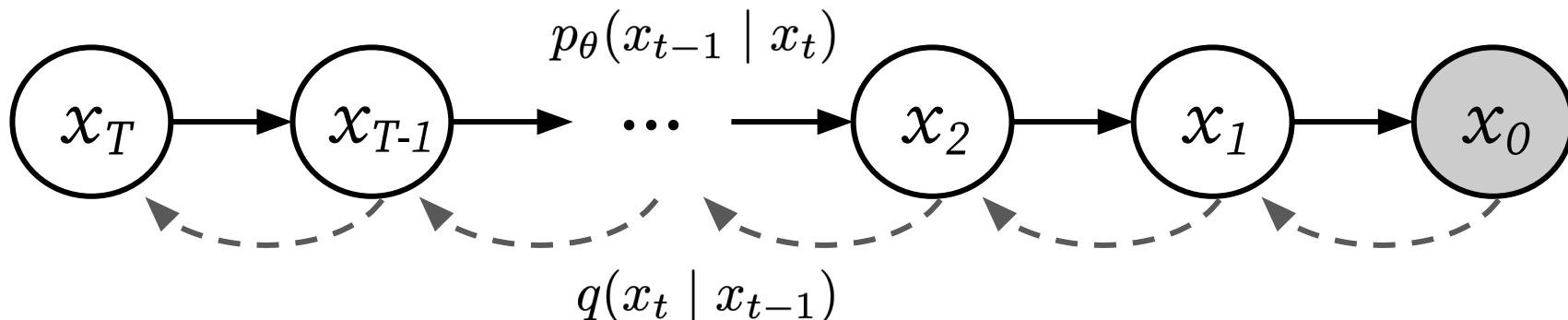
While the variational family of *posterior approximations* can be written:

$$q(x_1, \dots, x_T \mid x_0) = \prod_{t=1}^T q(x_t \mid x_{t-1})$$

(Another Markov chain!)

And sticking both together we get:

Similar to the “usual”
PGM diagram



Connections to Diffusion Models

Similar to VAE, you can:

- Write out the ELBO – has a particular form due to the structure of p and q .
- Then try to optimize it with respect to θ (and/or any VI parameters).

Connections to Diffusion Models

Before the DDPM paper, diffusion models as described here (roughly) already existed in prior work.

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 Mahowaratthan
Stanford University
Surya Ganguli
Stanford University

JASCHA@STANFORD.EDU
EWEISS@BERKELEY.EDU
NIKUM@STANFORD.EDU
SGANGULI@STANFORD.EDU

Abstract

A central problem in machine learning involves modeling complex data-sets using highly flexible families of probability distributions in which learning, sampling, inference, and evaluation are still analytically or computationally tractable. Historically, probabilistic models that are tractably learned usually achieve both flexibility and tractability. The essential idea, inspired by non-equilibrium statistical physics, is to perturb a model and slowly destroy structure in a data distribution through an iterative forward diffusion process. We then learn a reverse diffusion process that restores structure in the data, yielding a highly flexible and tractable generative model of the data. This approach allows us to rapidly learn sample from, and evaluate probabilities in, deep generative models with thousands of latent variable 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, these models are unable to apply descriptive structure in rich datasets. On the other hand, models that are *flexible* can be modeled to fit structure in arbitrary data. For example, we can define models in terms of any (non-negative) function $\phi(x)$ yielding the flexible distribution $p(x) = \frac{\phi(x)}{Z}$, where Z is a normalization constant. However, this normalization constant is generally intractable. Evaluating, training, or drawing samples from such flexible models typically requires a very expensive Monte Carlo process.

A variety of analytic techniques exist to make such models tractable, but at a cost. For instance, mean field theory and its expansions (T., 1982; Tanaka, 1998), variational Bayes (Jordan et al., 1999), contrastive divergence (Welling & Tipping, 2004), Langevin dynamics, minimum probability flow (Sohl-Dickstein et al., 2013(b)), minimum KL contraction (Lyu, 2011), proper scoring rules (Gneiting & Raftery, 2007; Parry et al., 2012), score matching (Welling & Teh, 2011), expectation-maximization (Eliezer & Rafferty, 2009), and many, many more. Non-parametric methods (Gershman & Blei, 2012) can also be very effective.

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.

Proceedings of the 32nd International Conference on Machine Learning, Lille, France, 2015. JMLR: W&CP volume 37. Copyright 2015 by the author(s).

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

Source: “Deep Unsupervised Learning using Nonequilibrium Thermodynamics”, Sohl-Dickstein et al., 2015

Connections to Diffusion Models

Before the DDPM paper, diffusion models as described here (roughly) already existed in prior work.

However, the DDPM paper does a couple of things:

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 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 best results are obtained by training on a weighted variational bound designed according to a novel connection between diffusion probabilistic models and variational autoencoders (VAEs). Our synthesis method naturally admits a progressive lossy decompression scheme that can be interpreted as a generalization of autoregressive decoding. On the unconditional CIFAR10 dataset, we obtain samples with PSNR 34.0 dB and SSIM 0.95 at 256x256. On 256x256 LSUN, we obtain sample quality similar to ProgressiveGAN. Our implementation is available at <https://github.com/jonathanho/diffusion>.

1 Introduction

Deep 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 autoencoders (VAEs) have synthesized striking image and audio samples [14, 27, 3, 38, 38, 25, 10, 32, 47, 20, 33, 45], and there have been notable 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.

Source: "Denoising Diffusion Probabilistic Models", Ho et al., 2019

Connections to Diffusion Models

Before the DDPM paper, diffusion models as described here (roughly) already existed in prior work.

However, the DDPM paper does a couple of things:

- (1) Shows that under a certain parameterization of the model, their algorithm is *equivalent to* denosing score matching.

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 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 best results are obtained by training on a weighted variational bound designed according to a novel connection between diffusion probabilistic models and denoising score matching. We show that diffusion probabilistic models naturally admit a progressive lossy decompression scheme that can be interpreted as a generalization of autoregressive decoding. On the unconditional CIFAR10 dataset, we obtain samples with PSNR = 33.45 dB at 256x256 resolution. On 256x256 LSUN, we obtain sample quality similar to ProgressiveGAN. Our implementation is available at <https://github.com/jonathanho/diffusion>.

1 Introduction

Deep 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 autoencoders (VAEs) have synthesized striking image and audio samples [14, 27, 3, 38, 38, 25, 10, 32, 47, 20, 33, 45], and there have been notable advances in energy-based modeling and score matching that have produced images comparable to those of GANs [11, 55].

The figure shows two sets of generated images. The left set, labeled 'Generated samples on CelebA-HQ 256 × 256', consists of two portrait photos of diverse individuals. The right set, labeled 'unconditional CIFAR10', is a grid of 16x16 smaller images showing various objects and scenes, such as animals, landscapes, and indoor settings.

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.

Source: "Denoising Diffusion Probabilistic Models", Ho et al., 2019

Connections to Diffusion Models

Before the DDPM paper, diffusion models as described here (roughly) already existed in prior work.

However, the DDPM paper does a couple of things:

- (1) Shows that under a certain parameterization of the model, their algorithm is *equivalent* to denosing score matching.
 - Shows optimizing the ELBO \Rightarrow equivalent to the “*sum of Fisher divergences*” loss from before!
 - And that sampling from the forward model \Rightarrow equivalent to Langevin dynamics on a learned score function!

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

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 non-equilibrium thermodynamics. Our best results are obtained by training on a weighted variational bound designed according to a novel connection between diffusion probabilistic models and denoising score matching. The resulting models are able to naturally admit a progressive lossy decompression scheme that can be interpreted as a generalization of autoregressive decoding. On the unconditional CIFAR10 dataset, we achieve a top-5 error of 4.9% at 256x256 resolution. On 256x256 LSUN, we obtain sample quality similar to ProgressiveGAN. Our implementation is available at <https://github.com/jonathanho/diffusion>.

1 Introduction

Deep 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 autoencoders (VAEs) have synthesized striking image and audio samples [14, 27, 3, 38, 38, 25, 10, 32, 47, 20, 33, 45], and there have been notable advances in energy-based modeling and score matching that have produced images comparable to those of GANs [11, 55].

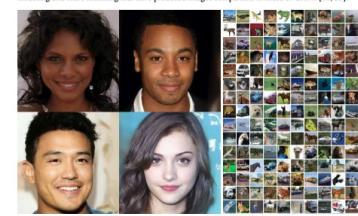


Figure 1: Generated samples on CelebA-HQ 256 × 256 (left) and unconditional CIFAR10 (right).

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

Source: “Denoising Diffusion Probabilistic Models”, Ho et al., 2019

Connections to Diffusion Models

Before the DDPM paper, diffusion models as described here (roughly) already existed in prior work.

However, the DDPM paper does a couple of things:

- (1) Shows that under a certain parameterization of the model, their algorithm is *equivalent to* denosing score matching.
- (2) Shows that under this parameterization (along with a few additional tricks to better handle image data), sample quality is can be very good.
(Possibly better than most/all previous generative models up until this point).

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

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 non-equilibrium thermodynamics. Our best results are obtained by training on a weighted variational bound designed according to a novel connection between diffusion probabilistic models and denoising score matching. This bound is also shown to be closely related to naturally admit a progressive lossy decompression scheme that can be interpreted as a generalization of autoregressive decoding. On the unconditional CIFAR10 dataset, we obtain samples with PSNR = 33.45 dB at 256x256 resolution. On 256x256 LSUN, we obtain sample quality similar to ProgressiveGAN. Our implementation is available at <https://github.com/jonathanho/diffusion>.

1 Introduction

Deep 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 autoencoders (VAEs) have synthesized striking image and audio samples [14, 27, 3, 38, 38, 25, 10, 32, 47, 20, 33, 45], and there have been notable advances in energy-based modeling and score matching that have produced images comparable to those of GANs [11, 55].

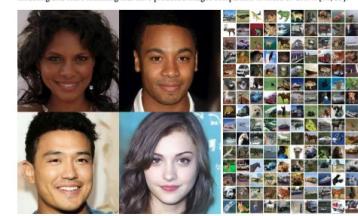


Figure 1: Generated samples on CelebA-HQ 256 × 256 (left) and unconditional CIFAR10 (right).
34th Conference on Neural Information Processing Systems (NeurIPS 2020), Vancouver, Canada.

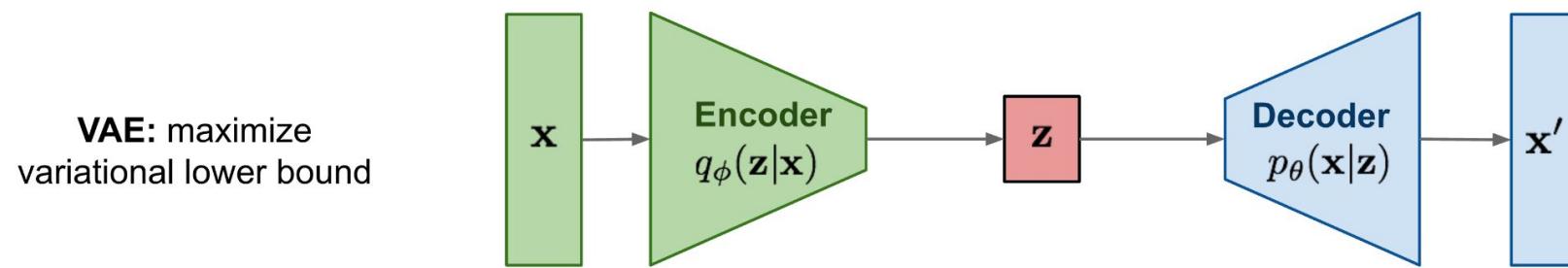
Source: "Denoising Diffusion Probabilistic Models", Ho et al., 2019

Connections to Diffusion Models

A comparison of generative modeling paradigms (from Lil'Log blog):

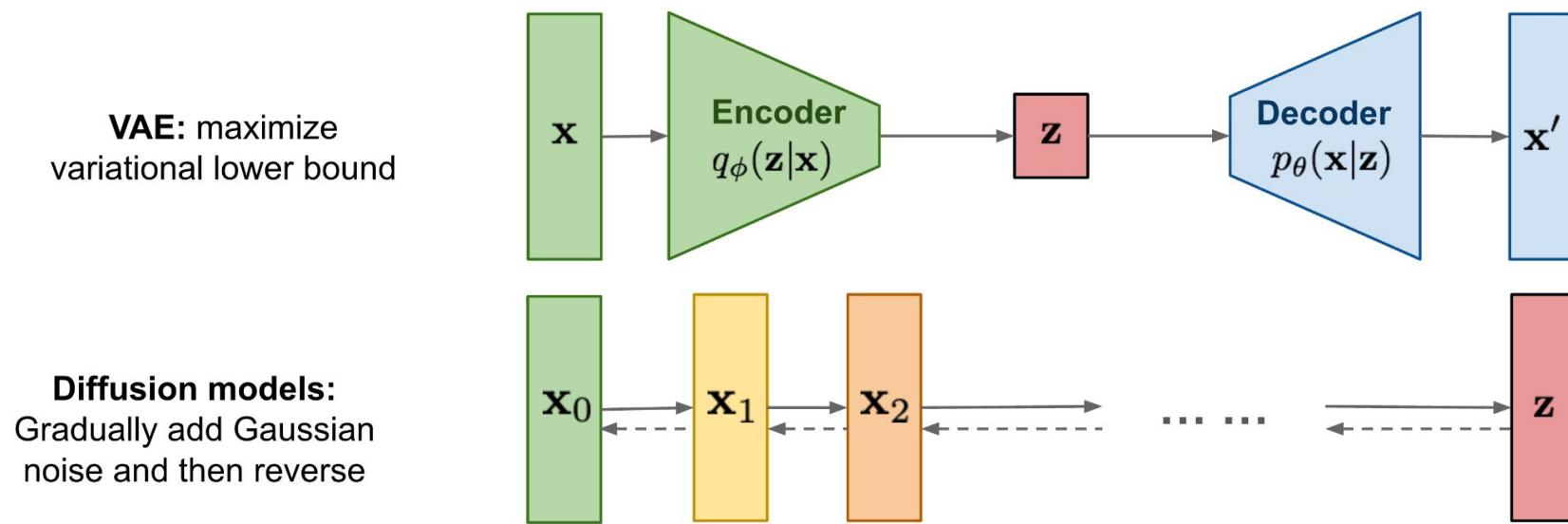
Connections to Diffusion Models

A comparison of generative modeling paradigms (from Lil'Log blog):



Connections to Diffusion Models

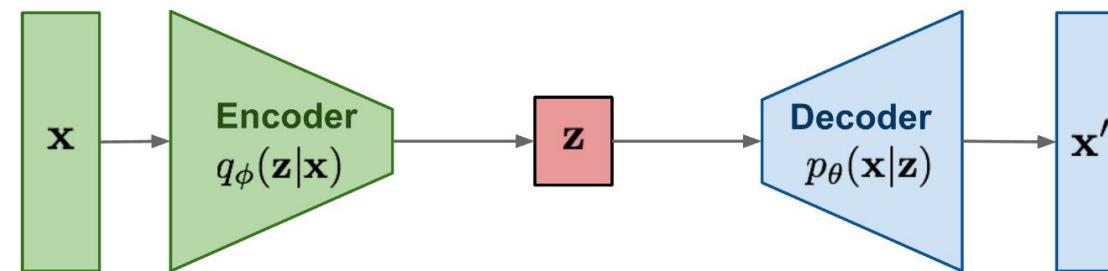
A comparison of generative modeling paradigms (from Lil'Log blog):



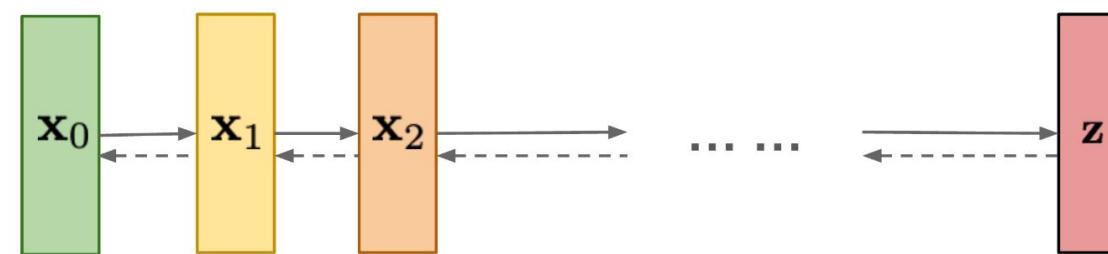
Connections to Diffusion Models

A comparison of generative modeling paradigms (from Lil'Log blog):

VAE: maximize variational lower bound



Diffusion models:
Gradually add Gaussian noise and then reverse



Flow-based models:
Invertible transform of distributions

