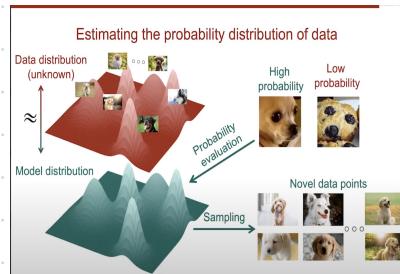


Diffusion and State-Based Generative Models

- A typical assumption in statistics and machine learning is that all data points in our training data set come from some underlying data distribution.



In other words, those data points are basically i.i.d. (independent and identically distributed) samples from given data distribution.

However, we don't have the analytical form of the data distribution of training dataset, and we need to estimate it.

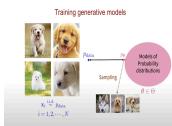
To estimate this data distribution, we need to create a model that represents parameterized probability distribution, which we call "the model distribution".

We hope to tune this model parameter to make sure this "model distribution" is close to the data distribution in a certain sense.

If this model distribution is very close to the data distribution, then we can use the model for many applications, such as generating unlimited number of novel data points just by sampling from this model distribution.

Because this model distribution provides a way to generate novel data points, we also refer to it as a "generative model".

Training Generative models



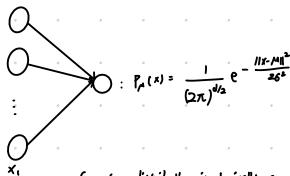
Our model provides a family of probability distributions. And we hope to find a single probability distribution inside huge family by minimizing the distances from P_{data} to P_{θ} .

Key challenge for building complex generative models.

① P_{data} is extremely complex for high dimensional data.

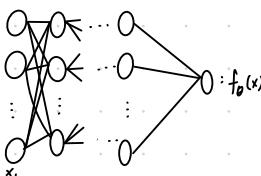
↳ How to build a complex model to fit the P_{data} ?

1. Gaussian distribution: It's too simple to approximate complicated P_{data} , but it serves as a good starting point.



Gaussian distribution is basically a computational graph that has two layers, where the first layer corresponds to the input data point and second layer is a single unit that basically gives you the probability density function of Gaussian distribution

2. But, Gaussian models are too simple. Let's leverage a bigger and deeper Network

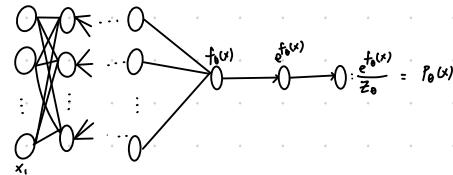


Deep Neural Network

But, it's actually non-trivial to use deep Neural Networks to directly represent a probability distribution. Because we typically view a deep Neural Network as a black box that converts a high dimensional input x to a typically one dimensional output $f_\theta(x)$.

The output value $f_\theta(x)$ does not directly model the distribution because it may not be positive everywhere. So, first step to convert $f_\theta(x)$ into a probability density is to ① take the exponential of the output, $e^{f_\theta(x)}$, so then the output becomes positive.

② Then, we normalize the output by dividing by a constant Z_0 , in order to construct a probability distribution which has positive values everywhere and also properly normalized.



Z_0 : Normalizing constant

$$Z_0 = \int e^{f_\theta(x)} dx$$

In the special case of Gaussian models

$$Z_\mu = \frac{1}{(2\pi)^{d/2}}$$

However in Deep Neural Network

$$\int e^{f_\theta(x)} dx \text{ becomes intractable to compute}$$

And, this difficulty is the unique challenge in deep generative modeling

Tackling the intractable normalizing constant

① Approximating the normalizing constant.

- Energy based models

Disadvantage: Inaccurate probability evaluation

② Using restricted Neural Network models.

- VAE

Disadvantages: Restricted model family

③ Modeling the Generation process Only

- GANs

Disadvantage: Cannot evaluate probabilities

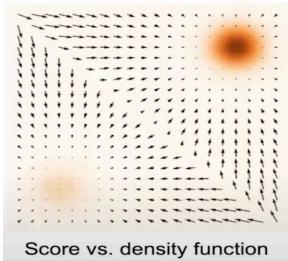
④ Working with score functions.

Working with Score functions

$P(x)$: Probability density function

$\nabla_x \log P(x)$: (Score) Score function

* Be careful this gradient is taken w.r.t. the r.v. 'x', not 'θ'



Score function is a vector field that gives the direction where the density function grows most quickly.

So, given the density function, we can compute the score function easily.

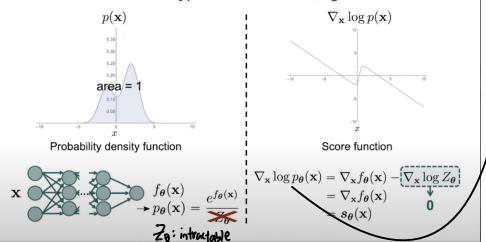
Conversely, with score function, we can recover pdf by computing integrals

Mathematically, this score function preserves all the information in the pdf.

Score functions bypass the normalizing constant

Flexible models

Score functions bypass the normalizing constant



$$\begin{aligned} \nabla_x \log p_\theta(x) &= \nabla_x \log \left(\frac{e^{f_\theta(x)}}{Z_\theta} \right) \\ &= \nabla_x \log(e^{f_\theta(x)}) - \nabla_x \log(Z_\theta) \\ &= \nabla_x f_\theta(x) - \nabla_x \log(Z_\theta) \\ &= \nabla_x f_\theta(x) - 0 \\ &= \nabla_x f_\theta(x) : \text{Score Model} \end{aligned}$$



We train $p_\theta(x)$ to estimate $p_{\text{data}}(x)$ (the underlying data distribution) using methods such as Maximum Likelihood.

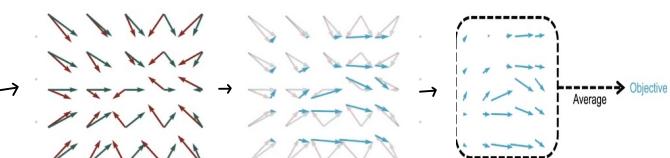
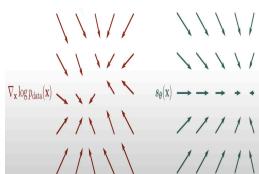
Score models can be estimated from data.

Given: $\{x_1, x_2, \dots, x_N\} \stackrel{\text{i.i.d.}}{\sim} p_{\text{data}}(x)$

Goal: $\nabla_x \log p_{\text{data}}(x)$

Score model: $S_\theta(x) : \mathbb{R}^d \rightarrow \mathbb{R}^d \approx \nabla_x \log p_{\text{data}}(x)$

Objective: How do we compare two vector fields of scores?



Let's recall that these two vector fields ($S_\theta(x)$ and $\nabla_x \log p_{\text{data}}(x)$) actually lie in the same space.

Mathematically, we can capture this intuition with the Fisher divergence objective (Fisher divergence is essentially an expected squared Euclidean distance between data score and the model score averaged over samples from the data)

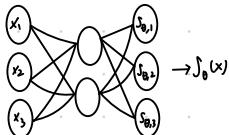
$$\frac{1}{2} \mathbb{E}_{p_{\text{data}}(x)} [\| \nabla_x \log p_{\text{data}}(x) - S_\theta(x) \|_2^2]$$

However, Fisher divergence cannot be directly computed because we don't have ground truth value of $p_{\text{data}}(x)$.

$$\text{Score Matching} : \mathbb{E}_{P_{\text{data}}(x)} \left[\frac{1}{2} \| J_\theta(x) \|^2_2 + \underbrace{\text{trace}(\nabla_x J_\theta(x))}_{\text{Jacobian of } J_\theta(x)} \right]$$

$$\begin{aligned} \frac{1}{2} \mathbb{E}_{P_{\text{data}}(x)} \left[\| \nabla_x \log P_{\text{data}}(x) - J_\theta(x) \|^2_2 \right] &= \mathbb{E}_{P_{\text{data}}(x)} \left[\frac{1}{2} \| J_\theta(x) \|^2_2 + \text{trace}(\nabla_x J_\theta(x)) \right] \\ &\approx \frac{1}{N} \sum_{i=1}^N \left[\frac{1}{2} \| J_\theta(x_i) \|^2_2 + \text{trace}(\nabla_x J_\theta(x_i)) \right] \end{aligned}$$

Deep Score Models



Compute $\text{trace}(\nabla_x J_\theta(x))$

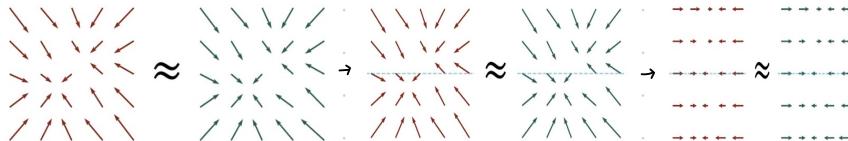
$$\nabla_x J_\theta(x) = \begin{pmatrix} \frac{\partial J_{\theta,1}(x)}{\partial x_1} & \frac{\partial J_{\theta,1}(x)}{\partial x_2} & \frac{\partial J_{\theta,1}(x)}{\partial x_3} \\ \frac{\partial J_{\theta,2}(x)}{\partial x_1} & \frac{\partial J_{\theta,2}(x)}{\partial x_2} & \frac{\partial J_{\theta,2}(x)}{\partial x_3} \\ \frac{\partial J_{\theta,3}(x)}{\partial x_1} & \frac{\partial J_{\theta,3}(x)}{\partial x_2} & \frac{\partial J_{\theta,3}(x)}{\partial x_3} \end{pmatrix} \quad \text{trace}(\nabla_x J_\theta(x)) = \begin{pmatrix} \frac{\partial J_{\theta,1}(x)}{\partial x_1} & \frac{\partial J_{\theta,1}(x)}{\partial x_2} & \frac{\partial J_{\theta,1}(x)}{\partial x_3} \\ \frac{\partial J_{\theta,2}(x)}{\partial x_1} & \frac{\partial J_{\theta,2}(x)}{\partial x_2} & \frac{\partial J_{\theta,2}(x)}{\partial x_3} \\ \frac{\partial J_{\theta,3}(x)}{\partial x_1} & \frac{\partial J_{\theta,3}(x)}{\partial x_2} & \frac{\partial J_{\theta,3}(x)}{\partial x_3} \end{pmatrix}$$

Computing $\text{trace}(\nabla_x J_\theta(x))$ requires the dimensionality of input data point. Hence, score matching in its naive form is not scalable.

Sliced Score Matching

Inuition: One dimensional problems should be easier.

Idea: Project onto random directions



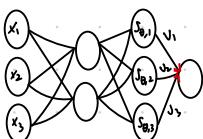
Random Objective: Sliced Fisher Divergence

$$: \frac{1}{2} \mathbb{E}_v \mathbb{E}_{P_{\text{data}}(x)} \left[(v^\top \nabla_x \log P_{\text{data}}(x) - v^\top J_\theta(x))^2 \right], \text{ where } v \text{ denotes the projection vector, } P_v \text{ denotes the distribution of projection directions}$$

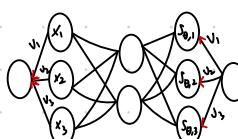
Sliced Score Matching:

$$\mathbb{E}_{P_v} \mathbb{E}_{P_{\text{data}}(x)} \left[v^\top \nabla_x J_\theta(x) + \frac{1}{2} (v^\top J_\theta(x))^2 \right]$$

Scalable



$$v^\top J_\theta(x)$$



$$v^\top \nabla_x (v^\top J_\theta(x))$$

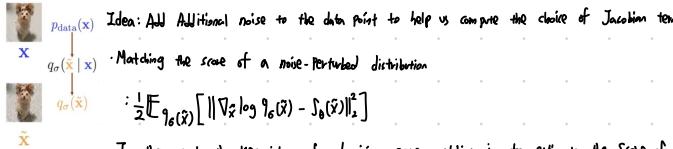
: One Backprop!

Sliced score matching

- ① Sample a minibatch of datapoints $\{x_1, x_2, \dots, x_N\} \sim p_{\text{data}}(x)$
 - ② Sample a minibatch of projection directions $\{v_1, v_2, \dots, v_N\} \sim p_v$
 - ③ Estimate the sliced score matching loss with Empirical Means:
- $$\frac{1}{N} \sum_{i=1}^N \left[v_i^\top \nabla_x s_\theta(x_i) v_i + \frac{\epsilon}{2} (\nabla_x s_\theta(x_i))^2 \right]$$
- ④ The projection distribution is typically Gaussian or Rademacher
 - ⑤ SGD
 - ⑥ Can use more projections per datapoint to boost performance

Denoising score matching

- Can also bypass the computational challenge of vanilla score matching.



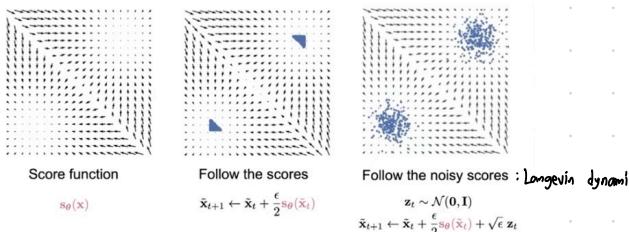
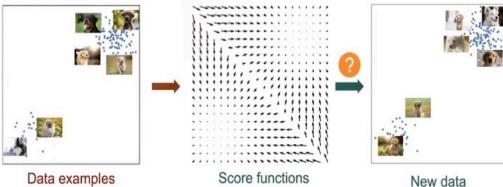
In other words, the key idea of denoising score matching is to estimate the score function of this noise data density ($S_\theta(\tilde{x})$) instead of the score function of the original data density ($S_\theta(x)$)

Denoising Score matching

$$\therefore \frac{1}{2} \mathbb{E}_{p_{\text{data}}(x)} \underbrace{\mathbb{E}_{q_\theta(\tilde{x}|x)} [\|\nabla_{\tilde{x}} \log q_\theta(\tilde{x}|x) - S_\theta(\tilde{x})\|_2^2]}_{\text{Scalable}}$$

- Cannot estimate scores of noise-free distributions.

Sampling from score functions. Langevin dynamics



Follow the scores: Move all points by following the directions predicted by the score function. However, it won't give us valid samples because all of those points will eventually collapse into each other.

Follow the noisy version of the score function: Inject Gaussian noise to score function and follow those noisy perturbed score function

Langevin dynamics sampling

Goal: Sample from $P(x)$ using only the score $\nabla_x \log P(x)$

Initialize $x^0 \sim \mathcal{U}(x)$

Repeat for $t \leftarrow 1, 2, \dots, T$

$$z^t \sim \mathcal{N}(0, I)$$

$$x^t \leftarrow x^{t-1} + \frac{\epsilon}{2} \nabla_x \log P(x^{t-1}) + \sqrt{\epsilon} z^t$$

If $\epsilon \rightarrow 0$ and $T \rightarrow \infty$, we are guaranteed to have $x^T \sim p(x)$

Langevin dynamics + Score estimation

$$S_\theta(x) \approx \nabla_x \log P(x)$$

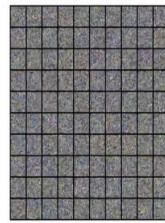
However, with this naive approach, we get following results.

Score matching + Langevin dynamics

CIFAR-10 data



Model samples



Why? Estimated scores are accurate in high data density regions. But for low data density regions, the estimated scores are not accurate at all.

Challenge in low data density regions



[Song and Ermon, NeurIPS 2019 (oral)]

It's because score matching compares the difference between the ground truth and the model only at samples from the data distribution.

So, in low data density regions, we don't have enough samples and therefore we don't have enough information to infer the score functions in those regions.

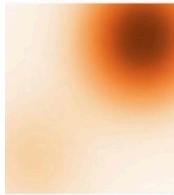
It's a huge obstacle for Langevin dynamics to provide high quality samples.

How can we address this challenge?

① Inject Gaussian noise to perturb our data points. After adding enough Gaussian noise, we perturb the data points to everywhere in the space. This means the size of low data density regions becomes smaller.

Improving score estimation by adding noise

Perturbed density



Perturbed scores



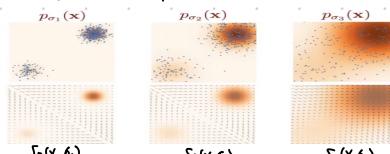
Estimated scores



[Song and Ermon, NeurIPS 2019 (oral)]

However, simply injecting Gaussian noise will not solve all the problems. Perturbed density no longer approximates the true data density.

: Using multiple sequence of different noise levels (learn low data density region efficiently)



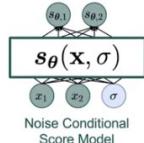
After obtaining noisy data sets, we want to estimate the underlying score function corresponding noisy data densities.

How can we estimate noisy score function?

Naive approach: train 3 networks, and each network is responsible for estimating the score function of a single noise level.

But, naive approach is not a scalable solution, because in practice, we might require much more noise levels.

Scalable solution: Noise Conditional Score Model



It's a simple modification to vanilla score model.

It takes noise level sigma as one additional input dimension to the model.

The output corresponds to the score function of the data density perturbed with noise level sigma

How to train this noise conditional score model? : We can leverage the idea of score matching

$$\frac{1}{N} \sum_{i=1}^N \lambda(\sigma_i) \mathbb{E}_{P_{\theta_i}(x)} \left[\nabla_x \log P_{\theta_i}(x) - S_{\theta}(x, \sigma_i) \right]_2^2$$

Weighting function Score matching loss

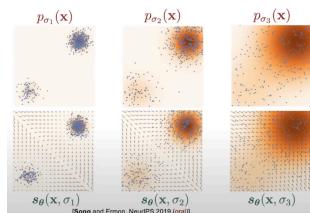
By minimizing this modified score matching loss, we will obtain accurate score estimation for all noise levels

Now, how to generate sample from noise-conditional Score model after training with the score matching loss?

: Annealed Langevin dynamics

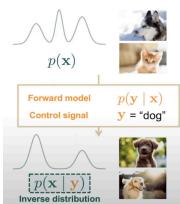
① First, apply Langevin dynamics to sample from the score model with the biggest perturbation noise, and the samples will be used as the initialization to sample from the score model of the next noise level

② Continue until it finally generate samples from the score function with the smallest noise level.



One remarkable property of a score-based generative model is the capability to control the generating process in a principled way

Control the generation process



How to obtain inverse distribution?

$$\text{Bayes' rule: } p(x|y) = \frac{p(x)p(y|x)}{p(y)}$$

Bayes rule for score functions:

$$\begin{aligned} \nabla_x \log p(x|y) &= \nabla_x \log p(x) + \nabla_x \log p(y|x) - \nabla_x \log p(y) \\ &= \nabla_x \log p(x) + \nabla_x \log p(y|x) - 0 \\ &\approx S_{\theta}(x) + \nabla_x \log p(y|x) \end{aligned}$$