

Generative Modelling

Goal :

Learn the distribution $p(x)$ over images x s.t

(1) generation : $x_{\text{new}} \sim p(x)$ is "good"

(2) density estimation : $p(x)$ should be high if x is similar to the data

(3) unsupervised representation learning : learn features that our datapoints have in common.

■ Question: how to represent $p(x)$?

Bayes Networks :

Recall :

$$\hookrightarrow (1) p(x_1, \dots, x_n) = p(x_1) p(x_2 | x_1) \dots p(x_n | x_1, \dots, x_{n-1}) \quad \begin{matrix} \text{Here } x_1, \dots, x_n \\ \text{could be pixels} \\ \text{of an image} \end{matrix}$$

Suppose $p(x_i)$ requires 1 parameter

$p(x_2 | x_1=0)$ requires 1 param, $p(x_2 | x_1=1)$ requires 1 param

Continuing this way, we require

$$1 + 2 + \dots + 2^{n-1} = 2^n - 1 \text{ parameters.}$$

This is a lot of parameters!

But if we impose a markov assumption...

$$\begin{aligned} p(x_1, \dots, x_n) &= p(x_1) p(x_2 | x_1) \dots p(x_n | x_1, \dots, x_{n-1}) \\ &= p(x_1) p(x_2 | x_1) p(x_3 | x_2) \dots p(x_n | x_{n-1}) \end{aligned}$$

This requires $2n-1$ parameters.

With this idea, we introduce Bayes Networks

for each random variable X_i , we learn $P(X_i | X_{A_i})$ where

X_{A_i} is a set of random variables.

Then,

$$p(x_1, \dots, x_n) = \prod_i p(x_i | X_{A_i})$$

We won't dig into the details any further.

■ Naive Bayes

Example: $Y=1$ implies the email is a spam, $Y=0$ otherwise
 (x_1, \dots, x_n) = an email.

Then $p(Y, x_1, \dots, x_n) = p(Y) \prod_{i=1}^n p(x_i | Y)$ by the Naive Bayes assumption.

We learn these parameters from the data:

$$p(Y=1 | x_1, \dots, x_n) = \frac{p(Y=1) \prod_{i=1}^n p(x_i | Y=1)}{\sum_{y \in \{0, 1\}} p(Y=y) \prod_{i=1}^n p(x_i | Y=y)}$$

Issue: The naive Bayes assumption is often too strong.



Here the Naive Bayes assumption is that

$$X_i \perp X_{-i} \mid Y$$

i.e. given Y , X_i is independent of X_{-i} , $j \in \{1, \dots, i\}$

Note : Discriminative vs Generative Models

Generative models: learn $P(Y, X)$ and $P(X)$

$$\text{Then, } P(Y|X) = \frac{P(Y, X)}{P(X)}$$

Discriminative models: Directly learn $P(Y|X)$ - no need to learn $P(X)$

Discriminative Models :

$$P(Y=1 | X; \theta) = f_\theta(x) \rightarrow \text{a parametrized function with parameters } \theta.$$

Often, we impose structure such as:

(1) linear dependence:

$$\text{define } z(\theta, x) = \theta_0 + \sum_{i=1}^n \theta_i x_i$$

$$\text{Then } P(Y=1 | X; \theta) = \sigma(z(\theta, x))$$

$$= \frac{1}{1 + \exp(-\theta_0 - \sum_{i=1}^n \theta_i x_i)}$$

Then, the decision boundary can be set to be $P(Y=1 | X; \theta) > 0.5$

(2) non-linear dependence:

Let $h(A, b, x) = f(Ax + b)$ be a non-linear transformation of the inputs

$$P_{\text{Neural}}(Y=1 | X; A, b, \theta) = f(\theta_0 + \sum_{i=1}^n \theta_i h(A, b, x_i))$$

These have more parameters to learn but are also more flexible.



Autoregressive Models

We order all the random variables like ordering the pixels from top-left to bottom-right.

Then,

$$P(x_1, \dots, x_n) = P(x_1) P(x_2 | x_1) \dots P(x_n | x_1, \dots, x_{n-1})$$

Maximum Likelihood Estimation

Assume that the data comes from P_{data} .

Dataset, D , of m samples from P_{data} .

Assume that the data instances are iid.

We have a family of models \mathcal{M} . We learn the optimal model $P_\theta \in \mathcal{M}$.

Goal is to learn P_{data} itself so we can answer any probabilistic inference query afterwards.

How do we do so?

$$(1) \text{ Minimize } D_{\text{KL}}(P \parallel Q) = \sum_x P(x) \log \frac{P(x)}{Q(x)}$$

Note: • $D_{\text{KL}}(P \parallel Q)$ is ≥ 0 and $= 0$ iff $P = Q$

This is because

$$D_{\text{KL}}(P \parallel Q) = \mathbb{E}_{x \sim P} \left[-\log \frac{Q(x)}{P(x)} \right] \geq -\log \mathbb{E}_{x \sim P} \left[\frac{Q(x)}{P(x)} \right] \quad (\text{by Jensen}) \\ = -\log (1) = 0$$

- Also, $D_{\text{KL}}(P \parallel Q) \neq D_{\text{KL}}(Q \parallel P)$

So, we $\boxed{\text{minimize } D_{\text{KL}}(P_{\text{data}} \parallel P_\theta)}$ which is 0 iff $P_{\text{data}} = P_\theta$.

$$\text{Note: } D_{\text{KL}}(P_{\text{data}} \parallel P_\theta) = \mathbb{E}_{x \sim P_{\text{data}}} \left[\log \left(\frac{P_{\text{data}}(x)}{P_\theta(x)} \right) \right] \\ = \underbrace{\mathbb{E}_{x \sim P_{\text{data}}} \left[\log(P_{\text{data}}(x)) \right]}_{\text{independent of } P_\theta} - \mathbb{E}_{x \sim P_{\text{data}}} \left[\log(P_\theta(x)) \right]$$

$$\therefore \arg \min_{P_\theta} D_{\text{KL}}(P_{\text{data}} \parallel P_\theta) = \arg \max_{P_\theta} \mathbb{E}_{x \sim P_{\text{data}}} \left[\log(P_\theta(x)) \right]$$

\uparrow
minimize KL divergence between P_{data} and P_θ

\uparrow
Maximize log probability of samples in the dataset!

This provides the justification for maximizing log likelihood in ML.

■ Issues: since we ignore $H(P_{\text{data}}) = -\mathbb{E}_{x \sim P_{\text{data}}} [\log P_{\text{data}}(x)]$,

we do not know how close we are to the optimum.

Also, we don't really know P_{data} , so how do we maximize log probability.

■ Solutions:

Option 1 : approximate $\mathbb{E}_{x \sim P_{\text{data}}} [\log(P_\theta(x))]$

$$\text{with } \mathbb{E}_D [\log(P_\theta(x))] = \frac{1}{|D|} \sum_{x \in D} \log(P_\theta(x))$$

Then, maximum likelihood learning is

$$\max_{\theta} \frac{1}{|D|} \sum_{x \in D} \log p_\theta(x)$$



for autoregressive models,
we maximize the
likelihood function

$$L(\theta, D) = \prod_{j=1}^m \prod_{i=1}^n p_{\text{neural}}(x_i^{(j)} | x_{<i}^{(j)}; \theta)$$

Aside: this is an example of

Monte Carlo Estimation
approximation

$$E_{x \sim p}[g(x)]$$

with

$$\frac{1}{T} \sum_{t=1}^T g(x^t) =: \hat{g}$$

where x^t are independent
samples from p .

→ this is an unbiased estimation:

$$E_p[\hat{g}] = E_p[g(x)]$$

→ convergence: by weak law of large no.

$$\hat{g} \rightarrow E_p[g(x)] \text{ as } T \rightarrow \infty.$$

→ variance:

$$V_p[\hat{g}] = V_p\left[\frac{1}{T} \sum_{i=1}^T g(x^i)\right] = \frac{V_p[g(x)]}{T}$$

so, $V_p[\hat{g}] \rightarrow 0$ as $T \rightarrow \infty$

Variational Autoencoders

10

Latent Variable Models

for our datapoints X , we aim to learn high-level features Z . If the features are learned properly, learning $P(X|Z)$ could be easier than learning $P(X)$ itself.

Example setup:

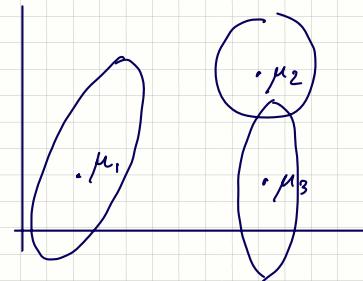
- (1) features $Z \sim \mathcal{N}(0, I)$. We hope that Z corresponds to meaningful ^{features}
- (2) Train $p_\theta(x|z) := \mathcal{N}(x | \mu_\theta(z), \Sigma_\theta(z))$
where μ_θ and Σ_θ are neural networks.

Now, we look at some examples of latent variable models.

(1) Mixture of Gaussians

$$Z \sim \text{Categorical}(1, \dots, k)$$

$$p_\theta(x|z=k) := \mathcal{N}(x | \mu_k, \Sigma_k)$$



Then, we can produce new data as

$$\begin{aligned} p(x) &= \sum_z p(x|z) \\ &= \sum_z p(z) p(x|z) \\ &= \sum_{i=1}^k p(z=i) \mathcal{N}(x | \mu_i, \Sigma_i) \end{aligned}$$

→ Alternative motivation is that each $\mathcal{N}(x | \mu_i, \Sigma_i)$ is a simple model and we are combining them.

→ Example construction:

$$\begin{aligned} \mu_\theta(z) &= \sigma(Az + c) \\ \Sigma_\theta(z) &= \text{diag}(\exp(\sigma(Bz + d))) = \begin{pmatrix} \exp(\sigma(b_1 z + d_1)) & 0 \\ 0 & \exp(\sigma(b_2 z + d_2)) \end{pmatrix} \end{aligned}$$

Issue: these are fairly difficult to learn

- We will use the following elements:
- (1) Evidence
 - (2) Reparameterization
 - (3) Amortization
- (2) Partially observed data: In most real-world settings, the datapoints X are observed in dataset D and the variables Z are not observed.
- If we learned $p_\theta(X, Z)$, we could do max. log likelihood:

$$\log \prod_{x \in D} p_\theta(x) = \sum_{x \in D} \log p_\theta(x) = \sum_{x \in D} \log \left(\sum_z p_\theta(x, z) \right)$$

But evaluating $\log \left(\sum_z p_\theta(x, z) \right)$ is intractable: if we have 30 binary features/latent, we need to sum over 2^{30} terms and if Z is continuous, then this is often intractable + gradient is hard to compute.
We have a few approaches then.

For computing $p_\theta(x)$:

Approach 1: Naive Monte Carlo

$$p_\theta(x) = \sum_z p_\theta(x, z) \approx |Z| \sum_{z \in Z} \frac{1}{|Z|} p_\theta(x, z) \\ = |Z| \mathbb{E}_{z \sim \text{Uniform}} \left[p_\theta(x, z) \right]$$

$$\text{Then, } \sum_z p_\theta(x, z) \approx |Z| \frac{1}{K} \sum_{j=1}^K p_\theta(x, z^{(j)})$$

Issue: for most z , $p_\theta(x, z)$ is very low

Approach 2: Importance sampling

$$p_\theta(x) = \sum_z p_\theta(x, z) = \sum_z \frac{q(z)}{q(z)} p_\theta(x, z) = \mathbb{E}_{z \sim q(z)} \left[\frac{p_\theta(x, z)}{q(z)} \right]$$

Then, we can do Monte Carlo to compute this:

$$p_\theta(x) \approx \frac{1}{K} \sum_{j=1}^K \frac{p_\theta(x, z^{(j)})}{q(z^{(j)})}$$

For computing $\log p_\theta(x)$:

How do we compute $\log p_\theta(x)$?

$$\log(p_\theta(x)) \approx \log \left(\frac{1}{K} \sum_{j=1}^K \frac{p_\theta(x, z^{(j)})}{q(z^{(j)})} \right)$$

$$\approx \log \left(\frac{p_\theta(x, z^{(j)})}{q(z^{(j)})} \right)$$

But this would be a bad approximation: as

$$\log \left(\mathbb{E}_{z \sim q(z)} \left[\frac{p_\theta(x, z)}{q(z)} \right] \right) \neq \underbrace{\mathbb{E}_{z \sim q(z)} \left[\log \left(\frac{p_\theta(x, z)}{q(z)} \right) \right]}$$

Approach 3: ELBO

by Jensen's inequality

$$\log p_\theta(x) = \underbrace{\log \left(\mathbb{E}_{z \sim q(z)} \left[\frac{p_\theta(x, z)}{q(z)} \right] \right)}_{\text{by importance sampling}} \geq \mathbb{E}_{z \sim q(z)} \left[\log \left(\frac{p_\theta(x, z)}{q(z)} \right) \right]$$

Evidence Lower Bound (ELBO)

→ Note:

$$\begin{aligned} \log p_\theta(x) &\geq \mathbb{E}_{z \sim q(z)} \left[\log \frac{p_\theta(x, z)}{q(z)} \right] \\ &= \mathbb{E}_{z \sim q(z)} \left[\log p_\theta(x, z) \right] - \underbrace{\mathbb{E}_{z \sim q(z)} \left[\log q(z) \right]}_{\text{Entropy of } q} \end{aligned}$$

→ Equality holds if $q = p_\theta(z|x)$

∴ Ideal choice of q is $q = p_\theta(z|x)$

Issue: what if $p_\theta(z|x)$ is hard to compute? How worse is this bound then?

→ Suppose $q(z)$ is any probability distribution.
Then,

$$D_{KL}(q(z) \parallel p_\theta(z|x)) = - \sum_z q(z) \log p_\theta(z|x) + \log p_\theta(x) - H(q) \geq 0$$

Thus, note that ELBO is equal to $\log p_\theta(x)$ if $q = p_\theta(z|x)$ (as $D_{KL}(q(z) \parallel p_\theta(z|x)) = 0$)

So, this KL-divergence term is a measure of the looseness of our bound.

→ we train $q_\phi(z)$

$$\text{Eg: } q_\phi(z) = \mathcal{N}(\phi_1, \phi_2)$$

Then, we do variational inference over ϕ so that $q_\phi(z)$ is close to $p_\theta(z|x)$.

Recap:

$$\log p_\theta(x) \geq \sum_z q_\phi(z) \log p_\theta(x, z) + H(q_\phi(z)) = \mathcal{L}_{\text{ELBO}}(x; \theta, \phi)$$

$$\text{and } \log p_\theta(x) = \mathcal{L}_{\text{ELBO}}(x; \theta, \phi) + D_{KL}(q_\phi(z) \parallel p_\theta(z|x))$$

Now, this was for one datapoint i.e. we made $q_\phi(z)$ close to $p_\theta(z|x)$. Then, for each datapoint $x^{(i)}$ we would have $\phi^{(i)}$ & train $q_{\phi^{(i)}}(z)$. If we did this, our MLE becomes

$$\ell(\theta; \mathcal{D}) = \sum_{x^{(i)} \in \mathcal{D}} \log p_\theta(x^{(i)}) \geq \sum_{x^{(i)} \in \mathcal{D}} \mathcal{L}(x^{(i)}; \theta, \phi^{(i)})$$

Issue : The problem is in optimizing this.

$$L(x^{(i)}, \theta, \phi^{(i)}) = \mathbb{E}_{z \sim q_{\phi^{(i)}}(z)} [\log p_{\theta}(x^{(i)}, z) - \log L_{\phi^{(i)}}(z)]$$

$$\begin{aligned} \nabla_{\phi} L(x^{(i)}, \theta, \phi^{(i)}) &= \nabla_{\phi} \mathbb{E}_{z \sim q_{\phi^{(i)}}(z)} [\log p_{\theta}(x^{(i)}, z)] \\ &= \mathbb{E}_{z \sim q_{\phi^{(i)}}(z)} [\nabla_{\phi} \log p_{\theta}(x^{(i)}, z)] \\ &\approx \frac{1}{k} \sum_k \nabla_{\phi} \log p_{\theta}(x^{(i)}, z^k) \end{aligned}$$

But

$\nabla_{\phi} L(x^{(i)}, \theta, \phi^{(i)})$ is hard as we have expectation wrt $q_{\phi^{(i)}}$.

Solution 1 : Reparametrization

In general suppose we want

$$\nabla_{\phi} \mathbb{E}_{z \sim q_{\phi}(z)} [r(z)] \text{ where } z \text{ is continuous.}$$

Now, note : if $q_{\phi}(z) = N(\mu, \sigma^2 I)$

$$\text{then } z = \mu + \sigma \varepsilon, \quad \varepsilon \sim N(0, I)$$

$$\Rightarrow z = g_{\phi}(\varepsilon) \text{ where } g \text{ is deterministic}$$

Then,

$$\begin{aligned} \mathbb{E}_{z \sim q_{\phi}(z)} [r(z)] &= \int q_{\phi}(z) r(z) dz \\ &= \mathbb{E}_{\varepsilon \sim N(0, I)} [r(g_{\phi}(\varepsilon))] \\ &= \int N(\varepsilon) r(\mu + \sigma \varepsilon) d\varepsilon \end{aligned}$$

$$\begin{aligned} \text{Then, } \nabla_{\phi} \mathbb{E}_{q_{\phi}(z)} [r(z)] &= \mathbb{E}_{\varepsilon} [\nabla_{\phi} r(g_{\phi}(\varepsilon))] \\ &\approx \frac{1}{k} \sum_k \nabla_{\phi} r(g_{\phi}(\varepsilon^k)) \end{aligned}$$

In our case, we assume $z = \mu + \sigma \varepsilon$ and then use this reparametrization technique.

Amortization

Instead of learning $\phi^{(i)}$ for $x^{(i)}$, $i \in \{1, \dots, 100\}$, we learn a single f_{λ} s.t

$$f_{\lambda}(x^{(i)}) \approx \phi^{(i)} \text{ for some good parameter } \phi^{(i)}.$$

Then, we approximate $q(z|x^{(i)})$ using $q_{\phi^{(i)}}(z)$ which is approximated with

$$q_{f_{\lambda}(x^{(i)})}(z) = q(z; f_{\lambda}(x^{(i)}))$$

Putting all of this together.....

Autoencoders

$$\begin{aligned}\text{Loss function : } \mathcal{L}(x; \theta, \phi) &= \mathbb{E}_{q_{\phi}(z|x)} [\log p_{\theta}(x|z) - \log q_{\phi}(z|x)] \\ &= \mathbb{E}_{q_{\phi}(z|x)} [\log p_{\theta}(x|z) - \log p(z) + \log p(z) - \log q_{\phi}(z|x)] \\ &= \mathbb{E}_{q_{\phi}(z|x)} [\log p_{\theta}(x|z) - \text{D}_{\text{KL}}(q_{\phi}(z|x) \parallel p(z))]\end{aligned}$$

Recipe : Take $x^{(i)}$, map it to $z^{(i)}$ using $\underbrace{q(z^{(i)}; f_{\lambda}(x^{(i)}))}_{\text{encoder}}$

Reconstruct $x^{(i)}$ by sampling $p_{\theta}(x|z^{(i)})$

Energy-based Models (EBMs)

When we try to parametrize $p_\theta(x)$, what are the constraints?

$$(1) p_\theta(x) \geq 0$$

$$(2) \sum_x p_\theta(x) = 1 \text{ or } \int p_\theta(x) dx = 1$$

The first is easy but generally $\sum_x p_\theta(x)$ need not be 1.
Solution:

$$\text{Set } p_\theta(x) = \frac{1}{Z(\theta)} g_\theta(x) = \frac{1}{\int g_\theta(x) dx} g_\theta(x)$$

Then, we choose g_θ s.t. we know $Z(\theta) = \int g_\theta(x) dx$

Example:

$$(1) g_{\mu, \sigma}(x) = \exp\left(-\frac{(x-\mu)^2}{2\sigma^2}\right)$$

$$Z(\mu, \sigma) = \sqrt{2\pi\sigma^2}$$

Gaussian

$$(2) g_\lambda(x) = e^{-\lambda x}$$

$$Z(\lambda) = \frac{1}{\lambda}$$

Exponential

(3) generally...

$$g_\theta(x) = h(x) \exp(\theta^\top T(x))$$

$$Z(\theta) = \exp(A(\theta)) = \int h(x) \exp(\theta^\top T(x)) dx$$

Exponential family

$$\text{Energy-based models: } p_\theta(x) = \frac{1}{\int \exp(f_\theta(x)) dx} \cdot \exp(f_\theta(x))$$

$$= \frac{1}{Z(\theta)} \exp(f_\theta(x))$$

called the partition function.

We use $\exp(f_\theta(x))$ to capture very large variations in probability + exponential family models have this structure.

Pros: flexibility in choosing $f_\theta(x)$

Cons:

- (1) Sampling difficulty

- (2) Evaluating is hard

- (3) No feature learning

- (4) Curse of dimensionality: computing $Z(\theta)$

→ however, some tasks don't require knowing $Z(\theta)$ like checking which sample is more likely:

$$\frac{p_\theta(x)}{p_\theta(x')} = \exp(f_\theta(x) - f_\theta(x'))$$

Goal : $\max \frac{\exp(f_\theta(x_{\text{train}}))}{Z(\theta)} = p_\theta(x_{\text{train}})$

Training Approach 1 : Monte Carlo estimate

instead of calculating $Z(\theta)$ exactly,
use contrastive divergence algorithm:

sample $x_{\text{sample}} \sim p_\theta$

Then, take step along

$$\nabla_\theta (f_\theta(x_{\text{train}}) - f_\theta(x_{\text{sample}}))$$

Make training data more likely
sample

This can be derived:

$$\begin{aligned} \bar{\nabla}_\theta \log p_\theta(x) &= \nabla_\theta f_\theta(x_{\text{train}}) - \nabla_\theta \log Z(\theta) \\ &= \nabla_\theta f_\theta(x_{\text{train}}) - \frac{\nabla_\theta Z(\theta)}{Z(\theta)} \\ &= \nabla_\theta f_\theta(x_{\text{train}}) - \frac{1}{Z(\theta)} \int \nabla_\theta \exp(f_\theta(x)) dx \\ &= \nabla_\theta f_\theta(x_{\text{train}}) - \frac{1}{Z(\theta)} \int (\nabla_\theta f_\theta(x)) e^{f_\theta(x)} dx \\ &= \nabla_\theta f_\theta(x_{\text{train}}) - \int p_\theta(x) \cdot \nabla_\theta f_\theta(x) dx \\ &= \nabla_\theta f_\theta(x_{\text{train}}) - \mathbb{E}_{x \sim p_\theta} [\nabla_\theta f_\theta(x)] \\ &\approx \nabla_\theta f_\theta(x_{\text{train}}) - \nabla_\theta f_\theta(x_{\text{sample}}) \end{aligned}$$

Sampling Approach

(1) Markov Chain Monte Carlo (MCMC):

Algorithm:

Initialize x^0 , $t = 0$

Let $x' = x^t + \text{noise}$:

if $f_\theta(x') > f_\theta(x^t)$, $x^{t+1} := x'$

else $x^{t+1} := x'$ with probability $\exp(f_\theta(x') - f_\theta(x^t))$

Issue: takes a long time to converge.

(2) Langevin MCMC:

Suppose we can compute the score function $\nabla_x \log p_\theta(x)$

Let $\pi(x)$ be a prior dist. that is easy to sample from.

Algorithm:

$x^0 \sim \pi(x)$

$x^{t+1} \sim x^t + \varepsilon \nabla_x \log p_\theta(x^t) + \sqrt{2\varepsilon} z^t$, $t = 0, \dots, T-1$

if $\varepsilon \rightarrow 0$ & $T \rightarrow \infty$, then $x^T \sim p_\theta(x)$.



$$\text{Note : } \nabla_x \log p_\theta(x) = \nabla_x f_\theta(x) - \underbrace{\nabla_x \log Z(\theta)}_0 \\ = \nabla_\theta f_\theta(x)$$

Score Matching

Recall : energy based models : $p_\theta(x) = \frac{\exp(f_\theta(x))}{Z(\theta)}$

Def : (Stein) score function

$$\begin{aligned}s_\theta(x) &:= \nabla_x \log p_\theta(x) \\ &= \nabla_x f_\theta(x) - \nabla_x \log Z(\theta) \\ &= \nabla_x f_\theta(x)\end{aligned}$$

Example :

(1) Gaussian distribution:

$$p_\theta(x) = \frac{1}{\sqrt{2\pi}\sigma} e^{-(x-\mu)^2/2\sigma^2}$$

$$\text{Then, } s_\theta(x) = -\frac{x-\mu}{\sigma^2}$$

(2) Gamma distribution:

$$p_\theta(x) = \frac{\beta^\alpha}{\Gamma(\alpha)} x^{\alpha-1} e^{-\beta x}$$

$$\text{Then, } s_\theta(x) = \frac{\alpha-1}{x} - \beta$$

Note :

(1) $s_\theta(x)$ is independent of $Z(\theta)$, the partition function.

Recall : fisher divergence between $p(x)$ and $q(x)$

$$D_F(p \parallel q) := \frac{1}{2} \mathbb{E}_{x \sim p} [\|\nabla_x \log p(x) - \nabla_x \log q(x)\|_2^2]$$

Score matching :

Minimize the fisher divergence between $p_{\text{data}}(x)$ and the EBM $p_\theta(x) \propto \exp(s_\theta(x))$:

$$\min. D_F(p_{\text{data}} \parallel p_\theta)$$

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

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

How do we compute this? Consider the univariate case for simplicity:

$$\begin{aligned}
 & \frac{1}{2} \mathbb{E}_{x \sim p_{\text{data}}} \left[(\nabla_x \log p_{\text{data}}(x) - \nabla_x \log p_\theta(x))^2 \right] \\
 &= \frac{1}{2} \int p_{\text{data}}(x) (\nabla_x \log p_{\text{data}}(x) - \nabla_x \log p_\theta(x))^2 dx \\
 &= \frac{1}{2} \int p_{\text{data}}(x) (\nabla_x \log p_{\text{data}}(x))^2 dx - \int p_{\text{data}}(x) \nabla_x \log p_{\text{data}}(x) \nabla_x \log p_\theta(x) dx \\
 &\quad + \frac{1}{2} \int p_{\text{data}}(x) (\nabla_x \log p_\theta(x))^2 dx
 \end{aligned}$$

$$\begin{aligned}
 \text{Now } & - \int p_{\text{data}}(x) \nabla_x \log p_{\text{data}}(x) \nabla_x \log p_\theta(x) dx \\
 &= - \int p_{\text{data}}(x) \frac{1}{p_{\text{data}}(x)} \nabla_x p_{\text{data}}(x) \nabla_x \log p_\theta(x) dx \\
 &= - p_{\text{data}}(x) \nabla_x \log p_\theta(x) \Big|_{-\infty}^{\infty} + \int p_{\text{data}}(x) \nabla_x^2 \log p_\theta(x) dx \\
 &\quad (\text{assuming } p_{\text{data}}(x) \rightarrow 0 \text{ as } x \rightarrow \pm\infty) \\
 &= \int p_{\text{data}}(x) \nabla_x^2 \log p_\theta(x) dx
 \end{aligned}$$

Then,

$$\begin{aligned}
 & \frac{1}{2} \mathbb{E}_{x \sim p_{\text{data}}} \left[(\nabla_x \log p_{\text{data}}(x) - \nabla_x \log p_\theta(x))^2 \right] \\
 &= \frac{1}{2} \underbrace{\int p_{\text{data}}(x) (\nabla_x \log p_{\text{data}}(x))^2 dx}_{\text{constant wrt } \theta} + \frac{1}{2} \int p_{\text{data}}(x) (\nabla_x \log p_\theta(x))^2 dx \\
 &\quad + \int p_{\text{data}}(x) \nabla_x^2 \log p_\theta(x) dx \\
 &= \mathbb{E}_{x \sim p_{\text{data}}} \left[\frac{1}{2} (\nabla_x \log p_\theta(x))^2 + \nabla_x^2 \log p_\theta(x) \right] + \text{constant}
 \end{aligned}$$

Similarly, in the multivariate case:

$$\begin{aligned}
 & \frac{1}{2} \mathbb{E}_{x \sim p_{\text{data}}} \left[\| \nabla_x \log p_{\text{data}}(x) - \nabla_x \log p_\theta(x) \|_2^2 \right] \\
 &= \mathbb{E}_{x \sim p_{\text{data}}} \left[\frac{1}{2} \| \nabla_x \log p_\theta(x) \|_2^2 + \text{tr}(\nabla_x^2 \log p_\theta(x)) \right] + \text{constant}
 \end{aligned}$$

Final objective: minimize:

$$\approx \frac{1}{n} \sum_{i=1}^n \left[\frac{1}{2} \| \nabla_{x_i} f_\theta(x_i) \|_2^2 + \text{tr}(\nabla_{x_i}^2 f_\theta(x_i)) \right]$$

Issue: $\nabla_x^2 \log p_\theta(x)$ is computationally expensive

Score-based Models

The ideas in score matching for EBMs can be generalized:

Paper : Generative Modelling by Estimating
Gradients of the data distribution
(Score-based Models)

Score-based Generative Modeling

Dataset : $\{x_i \in \mathbb{R}^D\}_{i=1}^n$
(iid samples)

Data Distribution : $p_{\text{data}}(x)$

Score of a probability density $p(x)$: $\nabla_x \log p(x)$
Score Network : $s_\theta : \mathbb{R}^D \rightarrow \mathbb{R}^D$ → Train s_θ to
approximate the score
of $p_{\text{data}}(x)$

Goal : Learn a model to generate new samples from $p_{\text{data}}(x)$
2 steps → ① Score matching
② Langevin dynamics

① Score matching for score estimation :

Train NN $s_\theta(x)$ to estimate $\nabla_x \log p(x)$ without

training a model
to estimate $p_{\text{data}}(x)$
first.

goal : $\arg \min_{\theta} \frac{1}{2} \mathbb{E}_{p_{\text{data}}(x)} [\| s_\theta(x) - \nabla_x \log p_{\text{data}}(x) \|_2^2]$ → ①

② Denoising score matching :

First perturbs data point x with a pre-specified noise distribution:
 $q_\sigma(\tilde{x}|x)$ → (noising x)

$$\tilde{x} = x + \epsilon, \epsilon \sim \text{N}(0, \sigma^2 I)$$

$$q_\sigma(\tilde{x}|x) := \frac{1}{(2\pi)^{d/2} \sigma^d} \exp\left(-\frac{1}{2\sigma^2} \|\tilde{x} - x\|^2\right)$$

$$q_\sigma(\tilde{x}) := \int q_\sigma(\tilde{x}|x) p_{\text{data}}(x) dx$$

If $q_\sigma(\tilde{x}) \approx p_{\text{data}}(\tilde{x})$ (given noise is small), then we write ①
as (using Theorem 1)

$$\arg \min_{\theta} \frac{1}{2} \mathbb{E}_{q_\sigma(\tilde{x}|x) p_{\text{data}}(x)} [\| s_\theta(\tilde{x}) - \nabla_{\tilde{x}} \log q_\sigma(\tilde{x}|x) \|_2^2]$$

Theorem 1 :

$$\mathbb{E}_{q_\theta(x, \tilde{x})} \left[\frac{1}{2} \|s_\theta(\tilde{x}) - \frac{\partial}{\partial \tilde{x}} \log q_\theta(\tilde{x}|x)\|^2 \right] = \mathbb{E}_{q_\theta(\tilde{x})} \left[\frac{1}{2} \|s_\theta(\tilde{x}) - \frac{\partial}{\partial \tilde{x}} \log q_\theta(\tilde{x})\|^2 \right]$$

Proof :

$$\frac{\partial}{\partial \tilde{x}} \log q_\theta(\tilde{x}|x) = \frac{1}{\sigma^2} (\tilde{x} - x)$$

Start from RHS

$$\mathbb{E}_{q_\theta(\tilde{x})} \left[\frac{1}{2} \|s_\theta(\tilde{x})\|^2 \right] - A(\theta) + C_2$$

where $C_2 := \mathbb{E}_{q_\theta(\tilde{x})} \left[\frac{1}{2} \left\| \frac{\partial \log q_\theta(\tilde{x})}{\partial \tilde{x}} \right\|^2 \right]$ independent of θ

$$\text{and } A(\theta) := \mathbb{E}_{q_\theta(\tilde{x})} \left[\left\langle s_\theta(\tilde{x}), \frac{\partial \log q_\theta(\tilde{x})}{\partial \tilde{x}} \right\rangle \right]$$

$$= \int_{\tilde{x}} q_\theta(\tilde{x}) \left\langle s_\theta(\tilde{x}), \frac{\partial \log q_\theta(\tilde{x})}{\partial \tilde{x}} \right\rangle d\tilde{x}$$

$$= \int_{\tilde{x}} q_\theta(\tilde{x}) \left\langle s_\theta(\tilde{x}), \frac{\frac{\partial}{\partial \tilde{x}} q_\theta(\tilde{x})}{q_\theta(\tilde{x})} \right\rangle d\tilde{x}$$

$$= \int_{\tilde{x}} \left\langle s_\theta(\tilde{x}), \frac{\partial}{\partial \tilde{x}} q_\theta(\tilde{x}) \right\rangle d\tilde{x}$$

$$= \int_{\tilde{x}} \left\langle s_\theta(\tilde{x}), \frac{\partial}{\partial \tilde{x}} \int_x q_\theta(x) \underbrace{q_\theta(\tilde{x}|x)}_{p_{\text{data}}(x)} dx \right\rangle d\tilde{x}$$

$$= \int_{\tilde{x}} \left\langle s_\theta(\tilde{x}), \int_x q_\theta(x) \frac{\partial q_\theta(\tilde{x}|x)}{\partial \tilde{x}} dx \right\rangle d\tilde{x}$$

$$= \int_{\tilde{x}} \left\langle s_\theta(\tilde{x}), \int_x q_\theta(x) q_\theta(\tilde{x}|x) \frac{\partial}{\partial \tilde{x}} \log q_\theta(\tilde{x}|x) dx \right\rangle d\tilde{x}$$

$$= \int_{\tilde{x}} \int_x q_\theta(x) q_\theta(\tilde{x}|x) \left\langle s_\theta(\tilde{x}), \frac{\partial}{\partial \tilde{x}} \log q_\theta(\tilde{x}|x) \right\rangle dx d\tilde{x}$$

$$= \int_{\tilde{x}} \int_x q_\theta(\tilde{x}, x) \left\langle s_\theta(\tilde{x}), \frac{\partial}{\partial \tilde{x}} \log q_\theta(\tilde{x}|x) \right\rangle dx d\tilde{x}$$

$$= \mathbb{E}_{q_\theta(\tilde{x}, x)} \left[\left\langle s_\theta(\tilde{x}), \frac{\partial}{\partial \tilde{x}} \log q_\theta(\tilde{x}|x) \right\rangle \right]$$

∴ RHS becomes

$$\mathbb{E}_{q_\theta(\tilde{x})} \left[\frac{1}{2} \|s_\theta(\tilde{x})\|^2 \right] - \mathbb{E}_{q_\theta(x, \tilde{x})} \left[\left\langle s_\theta(\tilde{x}), \frac{\partial}{\partial \tilde{x}} \log q_\theta(\tilde{x}|x) \right\rangle \right] + C_2$$

Now, from LHS, we can write

$$\mathbb{E}_{q_\theta(\tilde{x})} \left[\frac{1}{2} \|s_\theta(\tilde{x})\|^2 \right] - \mathbb{E}_{q_\theta(x, \tilde{x})} \left[\left\langle s_\theta(\tilde{x}), \frac{\partial}{\partial \tilde{x}} \log q_\theta(\tilde{x}|x) \right\rangle \right] + C_3$$

② Langevin Dynamics (LD)

LD produces samples from $p(x)$ using only the score function $\nabla_x \log p(x)$.

fix a step size $\epsilon > 0$ and $\tilde{x}_0 \sim \pi(x)$ with π being a prior distribution

LD method (recursively) :

$$\tilde{x}_t := \tilde{x}_{t-1} + \frac{\epsilon}{2} \underbrace{\nabla_x \log p(\tilde{x}_{t-1})}_{\approx s_\theta(\tilde{x}_{t-1})} + \sqrt{\epsilon} z_t, \quad z_t \sim \mathcal{N}(0, I)$$

Distribution of $\tilde{x}_T \approx p(x)$ when $\epsilon \rightarrow 0, T \rightarrow \infty$





Diffusion Models

DDPM

Paper : Denoising Diffusion Probabilistic Models.

Notation :

→ Data distribution : $x_0 \sim q(x_0)$

→ Learned parametrised distribution : $p_\theta(x)$

Forward Process (Noising)

$$q(x_{1:T} | x_0) := \prod_{t=1}^T q(x_t | x_{t-1})$$

where

$$q(x_t | x_{t-1}) := \mathcal{N}\left(x_t ; \sqrt{1-\beta_t} x_{t-1}, \beta_t I\right) \rightarrow (1)$$

We can actually sample x_t , given x_0 , in one step.

To do so, define

$$\alpha_t = 1 - \beta_t$$

$$\overline{\alpha}_t = \prod_{s=1}^t \alpha_s$$

$$\text{Let } \varepsilon_0, \dots, \varepsilon_{t-1} \sim \mathcal{N}(0, 1)$$

Then,

$$x_t = \sqrt{1-\beta_t} x_{t-1} + \sqrt{\beta_t} \varepsilon_{t-1} \quad (\text{from (1)})$$

$$= \sqrt{\alpha_t} x_{t-1} + \sqrt{1-\alpha_t} \varepsilon_{t-1}$$

$$= \sqrt{\alpha_t} \left(\sqrt{\alpha_{t-1}} x_{t-2} + \sqrt{1-\alpha_{t-1}} \varepsilon_{t-2} \right) + \sqrt{1-\alpha_t} \varepsilon_{t-1}$$

$$= \sqrt{\alpha_t \alpha_{t-1}} x_{t-2} + \sqrt{\alpha_t} \sqrt{1-\alpha_{t-1}} \varepsilon_{t-2} + \sqrt{1-\alpha_t} \varepsilon_{t-1}$$

$$\text{Now, } \sqrt{\alpha_t} \sqrt{1-\alpha_{t-1}} \varepsilon_{t-2} \sim \mathcal{N}(0, \alpha_t(1-\alpha_{t-1}))$$

$$\text{and } \sqrt{1-\alpha_t} \varepsilon_{t-1} \sim \mathcal{N}(0, 1-\alpha_t)$$

$$\therefore \sqrt{\alpha_t} \sqrt{1-\alpha_{t-1}} \varepsilon_{t-2} + \sqrt{1-\alpha_t} \varepsilon_{t-1} \sim \mathcal{N}(0, \alpha_t(1-\alpha_{t-1}) + (1-\alpha_t))$$

$$\therefore x_t = \sqrt{\alpha_t \alpha_{t-1}} x_{t-2} + \sqrt{1-\alpha_t \alpha_{t-1}} \varepsilon_{t-2} \rightarrow (2)$$

$$\therefore x_t \sim q(x_t | x_0) = \mathcal{N}\left(x_t ; \sqrt{\alpha_t} x_0, (1-\alpha_t) I\right) \rightarrow (3)$$

Backward Process (denoising)

The reverse process is also a Markov chain with learned Gaussian noise

Start by sampling $x_T \sim p(x_T) = \mathcal{N}(x_T; 0, I)$

then,

$$p_\theta(x_{0:T}) = p(x_T) \prod_{t=1}^T p_\theta(x_{t-1} | x_t)$$

$$\text{where } p_\theta(x_{t-1} | x_t) = \mathcal{N}(x_{t-1}; \mu_\theta(x_t; t), \Sigma_\theta(x_t, t))$$

↑ ↑
learned.

■ Version 1 of training loss function:

Maximize log probability:

$$\begin{aligned} \log p_\theta(x_0) &= \log \int_{x_{1:T}} p_\theta(x_{0:T}) dx_{1:T} \\ &= \log \int_{x_{1:T}} p_\theta(x_{0:T}) \frac{q(x_{1:T} | x_0)}{q(x_{1:T} | x_0)} dx_{1:T} \\ &= \log \mathbb{E}_{q(x_{1:T} | x_0)} \left[\frac{p_\theta(x_{0:T})}{q(x_{1:T} | x_0)} \right] \\ &\geq \mathbb{E}_{q(x_{1:T} | x_0)} \left[\log \frac{p_\theta(x_{0:T})}{q(x_{1:T} | x_0)} \right] \end{aligned}$$

Then, we minimize:

$$\begin{aligned} -\mathbb{E}_{x_0 \sim q(\cdot)} [\log p_\theta(x_0)] &\leq \mathbb{E}_{q(\cdot)} \left[-\log \frac{p_\theta(x_{0:T})}{q(x_{1:T} | x_0)} \right] \\ &= \mathbb{E}_q \left[-\log p_\theta(x_T) - \sum_{t \geq 1} \log \frac{p_\theta(x_{t-1} | x_t)}{q(x_t | x_{t-1})} \right] \end{aligned}$$

∴ loss function is:

$$L = \mathbb{E}_{q(x_{0:T})} \left[-\log p_\theta(x_T) - \sum_{t \geq 1} \log \frac{p_\theta(x_{t-1} | x_t)}{q(x_t | x_{t-1})} \right]$$

Version 2 (with variance reduction)

$$\begin{aligned}
 L &= \mathbb{E}_q \left[-\log p_\theta(x_T) - \sum_{t \geq 1} \log \frac{p_\theta(x_{t-1} | x_t)}{q(x_t | x_{t-1})} \right] \\
 &= \mathbb{E}_q \left[-\log p_\theta(x_T) - \sum_{t \geq 1} \log \frac{p_\theta(x_{t-1} | x_t)}{q(x_t | x_{t-1})} - \log \frac{p_\theta(x_0 | x_1)}{q(x_1 | x_0)} \right] \\
 &= \mathbb{E}_q \left[-\log p_\theta(x_T) - \sum_{t \geq 1} \log \frac{p_\theta(x_{t-1} | x_t)}{q(x_{t-1} | x_t, x_0)} \frac{q(x_{t-1} | x_t, x_0)}{q(x_t | x_0)} \right. \\
 &\quad \left. - \log \frac{p_\theta(x_0 | x_1)}{q(x_1 | x_0)} \right] \\
 &\quad (\text{Bayes Law})
 \end{aligned}$$

$$\begin{aligned}
 &= \mathbb{E}_{q(x_0:T)} \left[-\log \frac{p_\theta(x_T)}{q(x_T | x_0)} - \log p_\theta(x_0 | x_1) - \sum_{t \geq 1} \log \frac{p_\theta(x_{t-1} | x_t)}{q(x_{t-1} | x_t, x_0)} \right] \\
 &= \mathbb{E}_{q(x_0:T)} \left[D_{KL}(q(x_T | x_0) \| p_\theta(x_T)) - \log p_\theta(x_0 | x_1) \right. \\
 &\quad \left. + \sum_{t \geq 1} D_{KL}(q(x_{t-1} | x_t, x_0) \| p_\theta(x_{t-1} | x_0)) \right]
 \end{aligned}$$

Now,

$$q(x_{t-1} | x_t, x_0) = \mathcal{N}(x_{t-1}; \tilde{\mu}_t(x_t, x_0), \tilde{\beta}_t I)$$

$$\text{where } \tilde{\mu}_t(x_t, x_0) = \frac{\sqrt{\alpha_{t-1}}}{1 - \alpha_t} \beta_t x_0 + \frac{\sqrt{\alpha_{t-1}}(1 - \alpha_{t-1})}{1 - \alpha_t} x_t$$

$$\text{and } \tilde{\beta}_t := \frac{1 - \alpha_{t-1}}{1 - \alpha_t} \beta_t$$

Then,

$$\begin{aligned}
 L &= \mathbb{E}_{q(x_0:T)} \left[\overbrace{D_{KL}(q(x_T | x_0) \| p_\theta(x_T))}^{L_T} - \log p_\theta(x_0 | x_1) \right. \\
 &\quad \left. + \sum_{t \geq 1} \underbrace{D_{KL}(q(x_{t-1} | x_t, x_0) \| p_\theta(x_{t-1} | x_0))}_{L_{t-1}} \right]
 \end{aligned}$$



Version 3 [DDPM] :

Consider β_t to be constant for all t .

Then, L_t is independent of σ (recall: $\rho_\theta(x_T) = N(0, I)$)
so we ignore it.

Recall :
$$\left\{ \begin{array}{l} \rho_\theta(x_{t-1}|x_t) = N(x_{t-1}; \mu_\theta(x_t; t), \Sigma_\theta(x_t, t)) \\ q(x_{t-1}|x_t, x_0) = N(x_{t-1}; \tilde{\mu}_t(x_t, x_0), \tilde{\beta}_t I) \\ \text{where } \tilde{\mu}_t(x_t, x_0) = \frac{\sqrt{\alpha_{t-1}}}{1-\bar{\alpha}_t} \beta_t x_0 + \frac{\sqrt{\alpha_t}(1-\bar{\alpha}_{t-1})}{1-\bar{\alpha}_t} x_t \\ \text{and } \tilde{\beta}_t := \frac{1-\bar{\alpha}_{t-1}}{1-\bar{\alpha}_t} \beta_t \end{array} \right.$$

set $\Sigma_\theta(x_t, t) := \sigma_t^2 I$

and set $\sigma_t^2 = \beta_t$ or $\sigma_t^2 = \tilde{\beta}_t$ \rightarrow either works
 ↑
 Better for
 $x_T \sim N(0, I)$
 So, $\Sigma_\theta(x_t, t)$ is not learned.

Then, $\rho_\theta(x_{t-1}|x_t) = N(x_{t-1}; \mu_\theta(x_t, t), \sigma_t^2 I)$

So,
 $L_{t-1} = \mathbb{E}_{q} \left[\text{KL} \left(q(x_{t-1}|x_t, x_0) \| \rho_\theta(x_{t-1}|x_0) \right) \right]$
 $= \mathbb{E}_{q} \left[\frac{1}{2\sigma_t^2} \| \tilde{\mu}_t(x_t, x_0) - \mu_\theta(x_t, t) \|_2^2 \right] + C$
 ↑
 constant

Now, in the forward process, we know

$$x_t \sim q(x_t|x_0) = N(x_t; \sqrt{\alpha_t} x_0, (1-\bar{\alpha}_t) I)$$

So, $x_t = \sqrt{\alpha_t} x_0 + \sqrt{1-\bar{\alpha}_t} \varepsilon$

$$\text{So, } x_0 = \frac{x_t - \sqrt{1-\bar{\alpha}_t} \varepsilon}{\sqrt{\alpha_t}}$$

Then,

$$\begin{aligned} L_{t-1} - C &= \mathbb{E}_{\varepsilon, x_0 \sim q} \left[\frac{1}{2\sigma_t^2} \| \tilde{\mu}_t \left(x_t, \frac{x_t - \sqrt{1-\bar{\alpha}_t} \varepsilon}{\sqrt{\alpha_t}} \right) - \mu_\theta(x_t, t) \|_2^2 \right] \\ &= \mathbb{E}_{\varepsilon, x_0 \sim q} \left[\frac{1}{2\sigma_t^2} \| \frac{1}{\sqrt{\alpha_t}} \left(x_t - \frac{\beta_t}{\sqrt{1-\bar{\alpha}_t}} \varepsilon \right) - \mu_\theta(x_t, t) \|_2^2 \right] \end{aligned}$$

Since x_t is input to μ_θ , we set μ_θ to be:

$$\mu_\theta(x_t, t) = \frac{1}{\sqrt{\alpha_t}} \left(x_t - \frac{\beta_t}{\sqrt{1-\bar{\alpha}_t}} \varepsilon_\theta(x_t, t) \right)$$

where ε_θ is a function approximator to predict ε from x_t .

Then,

$$L_{t-1} = \mathbb{E}_{x_0, \varepsilon} \left[\frac{\rho_t^2}{2 \sigma_t^2 \alpha_t (1 - \bar{\alpha}_t)} \| \varepsilon - \varepsilon_\theta (\sqrt{\bar{\alpha}_t} x_0 + \sqrt{1 - \bar{\alpha}_t} \varepsilon, t) \|^2 \right]$$

What about L_0 ?

We set the first step of the reverse process to be an independent discrete decoder derived from the Gaussian $\mathcal{N}(x_0; \mu_\theta(x_{1:T}), \sigma_\theta^2 I)$:

$$p_\theta(x_0 | x_1) := \prod_{i=1}^d \int_{\delta_-(x_0^i)}^{\delta_+(x_0^i)} \mathcal{N}(x; \mu_\theta^i(x_{1:T}), \sigma_\theta^2) dx$$

$$\text{where } \delta_+(x) = \begin{cases} \infty & \text{if } x=1 \\ x + \frac{1}{2\pi\sigma} & \text{if } x < 1 \end{cases}$$

$$\delta_-(x) = \begin{cases} -\infty & \text{if } x=-1 \\ x - \frac{1}{2\pi\sigma} & \text{if } x > -1 \end{cases}$$

and $d = \text{data dimensionality}$

At the first step, we sample $\mu_\theta(x_{1:T})$ noisily.

Altogether, Simplified training objective becomes:

$$L(\theta) := \mathbb{E}_{t, x_0, \varepsilon} \left[\| \varepsilon - \varepsilon_\theta (\sqrt{\bar{\alpha}_t} x_0 + \sqrt{1 - \bar{\alpha}_t} \varepsilon, t) \|^2 \right]$$

$$t \sim \text{Uniform}\{1, \dots, T\}$$

Algorithm 1 Training

- 1: **repeat**
- 2: $x_0 \sim q(x_0)$
- 3: $t \sim \text{Uniform}\{1, \dots, T\}$
- 4: $\varepsilon \sim \mathcal{N}(\mathbf{0}, \mathbf{I})$
- 5: Take gradient descent step on

$$\nabla_\theta \| \varepsilon - \varepsilon_\theta (\sqrt{\bar{\alpha}_t} x_0 + \sqrt{1 - \bar{\alpha}_t} \varepsilon, t) \|^2$$
- 6: **until** converged

Algorithm 2 Sampling

- 1: $\mathbf{x}_T \sim \mathcal{N}(\mathbf{0}, \mathbf{I})$
- 2: **for** $t = T, \dots, 1$ **do**
- 3: $\mathbf{z} \sim \mathcal{N}(\mathbf{0}, \mathbf{I})$ if $t > 1$, else $\mathbf{z} = \mathbf{0}$
- 4: $\mathbf{x}_{t-1} = \frac{1}{\sqrt{\bar{\alpha}_t}} \left(\mathbf{x}_t - \frac{1 - \bar{\alpha}_t}{\sqrt{1 - \bar{\alpha}_t}} \varepsilon_\theta(\mathbf{x}_t, t) \right) + \sigma_t \mathbf{z}$
- 5: **end for**
- 6: **return** \mathbf{x}_0