

Review of “Denoising Diffusion Probabilistic Models” by Jonathan Ho, Ajay Jain, and Pieter Abbeel

Summary by Camden Kitowski

Contents

1 Full Citation	1
2 Paper Summary	1
3 Concept Review	1
3.1 Math Derivations	1
3.1.1 Forward Process	1
3.1.2 ELBO and Reverse Process	3
3.2 GitHub Repository Link	6
4 Math Terms and Definitions	6

1 Full Citation

Jonathan Ho, Ajay Jain, and Pieter Abbeel. “Denoising diffusion probabilistic models”. In: *Advances in neural information processing systems* 33 (2020), pp. 6840–6851

2 Paper Summary

- **Diffusion Probabilistic Model:** A model trained using variational inference using a parameterized Markov chain to generate high quality samples after a finite time.
- **Forward Process:** This process gradually adds noise to data samples until the data is a normal distribution of noise. The closed form of the forward process:

$$q(\mathbf{x}_t | \mathbf{x}_0) = \mathcal{N}(\mathbf{x}_t; \sqrt{\bar{\alpha}_t} \mathbf{x}_0, (1 - \bar{\alpha}_t) \mathbf{I}).$$

- **Reverse Process:** This process gradually removes noise to obtain a new data sample. The reverse process follows:

$$p_\theta(\mathbf{x}_{0:T}) := p(\mathbf{x}_T) \prod_{t=1}^T p_\theta(\mathbf{x}_{t-1} | \mathbf{x}_t), \quad p_\theta(\mathbf{x}_{t-1} | \mathbf{x}_t) := \mathcal{N}(\mathbf{x}_{t-1}; \boldsymbol{\mu}_\theta(\mathbf{x}_t, t), \boldsymbol{\Sigma}_\theta(\mathbf{x}_t, t)).$$

3 Concept Review

3.1 Math Derivations

3.1.1 Forward Process

The forward process, denoted with q , is a fixed Markov chain. The forward process gradually adds noise from a Gaussian distribution to the data over t time steps with a variance schedule from β_1, \dots, β_T . Let x_0 be the original image and let x_t be the final image with isotropic Gaussian noise.

$$q(\mathbf{x}_{1:T} \mid \mathbf{x}_0) := \prod_{t=1}^T q(\mathbf{x}_t \mid \mathbf{x}_{t-1}), \quad q(\mathbf{x}_t \mid \mathbf{x}_{t-1}) := \mathcal{N}(\mathbf{x}_t; \sqrt{1 - \beta_t} \mathbf{x}_{t-1}, \beta_t \mathbf{I}) \quad (0)$$

Let $\alpha_t = 1 - \beta_t$
Let $\bar{\alpha}_t = \prod_{i=1}^t \alpha_i$

$$q(x_t \mid x_{t-1}) = \mathcal{N}(x_t, \sqrt{1 - \beta_t} x_{t-1}, \beta_t \mathbf{I}) \quad (1)$$

$$= \sqrt{1 - \beta_t} x_{t-1} + \sqrt{\beta_t} \epsilon \quad (2)$$

$$= \sqrt{\alpha_t} x_{t-1} + \sqrt{1 - \alpha_t} \epsilon \quad (3)$$

$$= \sqrt{\alpha_t \alpha_{t-1}} x_{t-2} + \sqrt{1 - \alpha_t \alpha_{t-1}} \epsilon \quad (4)$$

$$= \sqrt{\alpha_t \alpha_{t-1} \alpha_{t-2}} x_{t-3} + \sqrt{1 - \alpha_t \alpha_{t-1} \alpha_{t-2}} \epsilon \quad (5)$$

$$= \sqrt{\alpha_t \alpha_{t-1} \dots \alpha_1 \alpha_0} x_0 + \sqrt{1 - \alpha_t \alpha_{t-1} \dots \alpha_1 \alpha_0} \epsilon \quad (6)$$

$$= \sqrt{\bar{\alpha}_t} x_0 + \sqrt{1 - \bar{\alpha}_t} \epsilon \quad (7)$$

$$q(\mathbf{x}_t \mid \mathbf{x}_0) = \mathcal{N}(\mathbf{x}_t; \sqrt{\bar{\alpha}_t} \mathbf{x}_0, (1 - \bar{\alpha}_t) \mathbf{I}) \quad (8)$$

The following derivations break down steps 3 and 4.

$$\begin{aligned} x_t &= \sqrt{1 - \beta_t} x_{t-1} + \sqrt{\beta_t} \varepsilon_{t-1} \\ &= \sqrt{\alpha_t} x_{t-1} + \sqrt{1 - \alpha_t} \varepsilon_{t-1} \\ &= \sqrt{\alpha_t} \left(\sqrt{\alpha_{t-1}} x_{t-2} + \sqrt{1 - \alpha_{t-1}} \varepsilon_{t-2} \right) + \sqrt{1 - \alpha_t} \varepsilon_{t-1} \\ &= \sqrt{\alpha_t \alpha_{t-1}} x_{t-2} + \sqrt{\alpha_t (1 - \alpha_{t-1})} \varepsilon_{t-2} + \sqrt{1 - \alpha_t} \varepsilon_{t-1} \end{aligned}$$

Let $X = \sqrt{\alpha_t (1 - \alpha_{t-1})} \varepsilon_{t-2}$
Let $Y = \sqrt{1 - \alpha_t} \varepsilon_{t-1}$

Using reparameterization trick

$$\begin{aligned} 0 + \sqrt{\alpha_t (1 - \alpha_{t-1})} \varepsilon_{t-2} \dots \dots \dots \rightarrow X &\sim \mathcal{N}(0, \alpha_t (1 - \alpha_{t-1}) I) \\ 0 + \sqrt{1 - \alpha_t} \varepsilon_{t-1} \dots \dots \dots \rightarrow Y &\sim \mathcal{N}(0, (1 - \alpha_t) I) \end{aligned}$$

$$\text{If } X \sim \mathcal{N}(\mu_X, \sigma_X^2) \quad Y \sim \mathcal{N}(\mu_Y, \sigma_Y^2)$$

$$Z = X + Y$$

$$\text{Then } Z \sim \mathcal{N}(\mu_X + \mu_Y, \sigma_X^2 + \sigma_Y^2)$$

$$\begin{aligned} Z \sim \mathcal{N}(0, \sigma_X^2 + \sigma_Y^2) \dots \dots \dots \rightarrow \sigma_X^2 + \sigma_Y^2 &= \alpha_t (1 - \alpha_{t-1}) + (1 - \alpha_t) \\ &= \alpha_t - \alpha_t \alpha_{t-1} + 1 - \alpha_t \end{aligned}$$

$$Z \sim \mathcal{N}(0, (1 - \alpha_t \alpha_{t-1}) I) \leftarrow \dots \dots \dots = 1 - \alpha_t \alpha_{t-1}$$

$$\begin{aligned}
&= \sqrt{\alpha_t \alpha_{t-1}} x_{t-2} + \sqrt{1 - \alpha_t \alpha_{t-1}} \bar{\varepsilon}_{t-2} \\
&\vdots \\
&= \sqrt{\alpha_t \alpha_{t-1} \cdots \alpha_1} x_0 + \sqrt{1 - \alpha_t \alpha_{t-1} \cdots \alpha_1} \varepsilon \\
&= \sqrt{\bar{\alpha}_t} x_0 + \sqrt{1 - \bar{\alpha}_t} \varepsilon \quad \text{same result in line 7 above}
\end{aligned}$$

3.1.2 ELBO and Reverse Process

The loss function is derived below. The loss function is the $-\log(p_\theta(x_0))$. However, this is intractable because it depends on all other time steps before x_0 starting at x_T . This would mean keeping track of $t-1$ random variables, which is not possible in practice. The solution is to compute the variational lower bound for this objective. The overall setup of diffusion models is close to the setup of variational autoencoders.

$$\begin{aligned}
-\log(p_\theta(x_0)) &\leq -\log(p_\theta(x_0)) + \frac{D_{KL}(q(x_{1:T} | x_0) || p_\theta(x_{1:T} | x_0))}{\log\left(\frac{q(x_{1:T} | x_0)}{p_\theta(x_{1:T} | x_0)}\right)} \\
&\quad \downarrow \\
&\quad \log\left(\frac{q(x_{1:T} | x_0)}{p_\theta(x_{1:T} | x_0)}\right) \\
&\quad p_\theta(x_{1:T} | x_0) \\
&\quad \frac{p_\theta(x_0 | x_{1:T}) p_\theta(x_{1:T})}{p_\theta(x_0)} \quad \text{Bayes' rule} \\
&\quad \frac{p_\theta(x_0, x_{1:T})}{p_\theta(x_0)} \quad \text{Joint probability} \\
&\quad \frac{p_\theta(x_{0:T})}{p_\theta(x_0)} \quad \text{Combine probability} \\
&\quad \log\left(\frac{q(x_{1:T} | x_0)}{\frac{p_\theta(x_{0:T})}{p_\theta(x_0)}}\right) \rightarrow \log\left(\frac{q(x_{1:T} | x_0)}{p_\theta(x_{0:T})}\right) + \log(p_\theta(x_0)) \\
-\log(p_\theta(x_0)) &\leq -\log(p_\theta(x_0)) + \log\left(\frac{q(x_{1:T} | x_0)}{p_\theta(x_{0:T})}\right) + \log(p_\theta(x_0)) \\
-\log(p_\theta(x_0)) &\leq \log\left(\frac{q(x_{1:T} | x_0)}{p_\theta(x_{0:T})}\right)
\end{aligned}$$

$$\mathbb{E}[-\log p_\theta(x_0)] \leq \mathbb{E}_q\left[-\log \frac{p_\theta(x_{0:T})}{q(x_{1:T} | x_0)}\right] =: L$$

ELBO Derivation Continued

$$\begin{aligned}
L &= \mathbb{E}_q \left[\log \left(\frac{q(x_{1:T} | x_0)}{p_\theta(x_{0:T})} \right) \right] \\
&= \mathbb{E}_q \left[\log \left(\frac{\prod_{t=1}^T q(x_t | x_{t-1})}{p(x_T) \prod_{t=1}^T p_\theta(x_{t-1} | x_t)} \right) \right] && \text{Product form} \\
&= \mathbb{E}_q \left[-\log(p(x_T)) + \sum_{t=1}^T \log \left(\frac{q(x_t | x_{t-1})}{p_\theta(x_{t-1} | x_t)} \right) \right] && \text{Turn products into sums} \\
&= \mathbb{E}_q \left[-\log(p(x_T)) + \sum_{t=2}^T \log \left(\frac{q(x_t | x_{t-1})}{p_\theta(x_{t-1} | x_t)} \right) + \log \left(\frac{q(x_1 | x_0)}{p_\theta(x_0 | x_1)} \right) \right] && \text{Split up summation} \\
&= \mathbb{E}_q \left[-\log(p(x_T)) + \sum_{t=2}^T \log \left(\frac{q(x_{t-1} | x_t, x_0) q(x_t | x_0)}{p_\theta(x_{t-1} | x_t) q(x_{t-1} | x_0)} \right) + \log \left(\frac{q(x_1 | x_0)}{p_\theta(x_0 | x_1)} \right) \right] && \text{Bayes' rule and condition on } x_0 \\
&= \mathbb{E}_q \left[-\log(p(x_T)) + \sum_{t=2}^T \log \left(\frac{q(x_{t-1} | x_t, x_0)}{p_\theta(x_{t-1} | x_t)} \right) + \sum_{t=2}^T \log \left(\frac{q(x_t | x_0)}{q(x_{t-1} | x_0)} \right) + \log \left(\frac{q(x_1 | x_0)}{p_\theta(x_0 | x_1)} \right) \right] && \text{Split up summation} \\
&= \mathbb{E}_q \left[-\log(p(x_T)) + \sum_{t=2}^T \log \left(\frac{q(x_{t-1} | x_t, x_0)}{p_\theta(x_{t-1} | x_t)} \right) + \log \left(\frac{q(x_T | x_0)}{q(x_1 | x_0)} \right) + \log \left(\frac{q(x_1 | x_0)}{p_\theta(x_0 | x_1)} \right) \right] && \text{Simplify summation} \\
&= \mathbb{E}_q \left[-\log(p(x_T)) + \sum_{t=2}^T \log \left(\frac{q(x_{t-1} | x_t, x_0)}{p_\theta(x_{t-1} | x_t)} \right) + \log q(x_T | x_0) - \log p_\theta(x_0 | x_1) \right] && \text{Log rules and cancel} \\
&= \mathbb{E}_q \left[\log \left(\frac{q(x_T | x_0)}{p(x_T)} \right) + \sum_{t=2}^T \log \left(\frac{q(x_{t-1} | x_t, x_0)}{p_\theta(x_{t-1} | x_t)} \right) - \log p_\theta(x_0 | x_1) \right] && \text{Combine terms} \\
&= \mathbb{E}_q \left[\underbrace{D_{\text{KL}}(q(\mathbf{x}_T | \mathbf{x}_0) \| p(\mathbf{x}_T))}_{L_T} + \sum_{t=2}^T \underbrace{D_{\text{KL}}(q(\mathbf{x}_{t-1} | \mathbf{x}_t, \mathbf{x}_0) \| p_\theta(\mathbf{x}_{t-1} | \mathbf{x}_t))}_{L_{t-1}} - \underbrace{\log p_\theta(\mathbf{x}_0 | \mathbf{x}_1)}_{L_0} \right] && \text{KL divergence notation}
\end{aligned}$$

L_T - Constant Term - Since q does not have any learnable parameters, this term will be constant during training. This term can be ignored.

L_{t-1} - Denoising Term - Reverse Process - The reverse process starts with an isotropic Gaussian image, x_t , and denoises the image to make a sample that resembles an image in the training data, x_0 . The reverse process is an iterative process to denoise the image with each new image dependent on the last image within the sequence. This term compares the actual noise (q) and predicted noise (p).

$$D_{\text{KL}}(q(\mathbf{x}_{t-1} | \mathbf{x}_t, \mathbf{x}_0) \| p_\theta(\mathbf{x}_{t-1} | \mathbf{x}_t))$$

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

$$\text{where} \quad \tilde{\mu}_t(x_t, x_0) := \frac{\sqrt{\bar{\alpha}_{t-1}}\beta_t}{1 - \bar{\alpha}_t}x_0 + \frac{\sqrt{\bar{\alpha}_t}(1 - \bar{\alpha}_{t-1})}{1 - \bar{\alpha}_t}x_t \quad \text{and} \quad \tilde{\beta}_t := \frac{1 - \bar{\alpha}_{t-1}}{1 - \bar{\alpha}_t}\beta_t, \quad x_0 = \frac{1}{\sqrt{\bar{\alpha}_t}}(x_t - \sqrt{1 - \bar{\alpha}_t}\varepsilon_t)$$

$$\begin{aligned}
\tilde{\mu}_t(x_t, x_0) &= \frac{\sqrt{\alpha_t}(1 - \bar{\alpha}_{t-1})}{1 - \bar{\alpha}_t} x_t + \frac{\sqrt{\bar{\alpha}_{t-1}}\beta_t}{1 - \bar{\alpha}_t} x_0 \\
\tilde{\mu}_t(x_t) &= \frac{\sqrt{\alpha_t}(1 - \bar{\alpha}_{t-1})}{1 - \bar{\alpha}_t} x_t + \frac{\sqrt{\bar{\alpha}_{t-1}}\beta_t}{1 - \bar{\alpha}_t} \cdot \frac{1}{\sqrt{\bar{\alpha}_t}} (x_t - \sqrt{1 - \bar{\alpha}_t}\varepsilon_t) \quad \text{Substitute } x_0 \\
&= \frac{\sqrt{\alpha_t}(1 - \bar{\alpha}_{t-1})}{1 - \bar{\alpha}_t} x_t + \frac{\sqrt{\bar{\alpha}_{t-1}}\beta_t}{1 - \bar{\alpha}_t} \cdot \frac{1}{\sqrt{\bar{\alpha}_{t-1}} \cdot \alpha_t} (x_t - \sqrt{1 - \bar{\alpha}_t}\varepsilon_t) \quad \text{Break up } \alpha_t \\
&= \frac{\sqrt{\alpha_t}(1 - \bar{\alpha}_{t-1})}{1 - \bar{\alpha}_t} x_t + \frac{\beta_t}{\sqrt{\alpha_t}(1 - \bar{\alpha}_t)} (x_t - \sqrt{1 - \bar{\alpha}_t}\varepsilon_t) \quad \text{Cancel terms} \\
&= \frac{\sqrt{\alpha_t}(1 - \bar{\alpha}_{t-1})}{1 - \bar{\alpha}_t} x_t + \frac{\beta_t}{\sqrt{\alpha_t}(1 - \bar{\alpha}_t)} x_t - \frac{\beta_t}{\sqrt{\alpha_t}(1 - \bar{\alpha}_t)} \sqrt{1 - \bar{\alpha}_t}\varepsilon_t \quad \text{Break up terms} \\
&= \frac{\sqrt{\alpha_t}(1 - \bar{\alpha}_{t-1})}{1 - \bar{\alpha}_t} x_t + \frac{1 - \alpha_t}{\sqrt{\alpha_t}(1 - \bar{\alpha}_t)} x_t - \frac{1 - \alpha_t}{\sqrt{\alpha_t}(1 - \bar{\alpha}_t)} \sqrt{1 - \bar{\alpha}_t}\varepsilon_t \quad \text{Substitute } \beta_t \\
&= \frac{\alpha_t}{\sqrt{\alpha_t}} \cdot \frac{(1 - \bar{\alpha}_{t-1})}{(1 - \bar{\alpha}_t)} x_t + \frac{1 - \alpha_t}{\sqrt{\alpha_t}(1 - \bar{\alpha}_t)} x_t - \frac{1 - \alpha_t}{\sqrt{\alpha_t}} \cdot \frac{\sqrt{1 - \bar{\alpha}_t}}{(1 - \bar{\alpha}_t)} \cdot \varepsilon_t \quad \text{Reformulate } \sqrt{\alpha_t} \\
&= \frac{\alpha_t - \bar{\alpha}_{t-1} \cdot \alpha_t}{\sqrt{\alpha_t}(1 - \bar{\alpha}_t)} x_t + \frac{1 - \alpha_t}{\sqrt{\alpha_t}(1 - \bar{\alpha}_t)} x_t - \frac{1 - \alpha_t}{\sqrt{\alpha_t}} \cdot \frac{1}{\sqrt{1 - \bar{\alpha}_t}} \cdot \varepsilon_t \quad \text{Reformulate } \frac{\sqrt{1 - \bar{\alpha}_t}}{(1 - \bar{\alpha}_t)} \\
&= \frac{\alpha_t - \bar{\alpha}_t + 1 - \alpha_t}{\sqrt{\alpha_t}(1 - \bar{\alpha}_t)} x_t - \frac{1}{\sqrt{\alpha_t}} \cdot \frac{1 - \alpha_t}{\sqrt{1 - \bar{\alpha}_t}} \cdot \varepsilon_t \quad \text{Combine terms} \\
&= \frac{1 - \bar{\alpha}_t}{\sqrt{\alpha_t}(1 - \bar{\alpha}_t)} x_t - \frac{1}{\sqrt{\alpha_t}} \cdot \frac{1 - \alpha_t}{\sqrt{1 - \bar{\alpha}_t}} \cdot \varepsilon_t \quad \text{Cancel terms} \\
&= \frac{1}{\sqrt{\alpha_t}} x_t - \frac{1}{\sqrt{\alpha_t}} \cdot \frac{1 - \alpha_t}{\sqrt{1 - \bar{\alpha}_t}} \cdot \varepsilon_t \quad \text{Cancel terms} \\
\tilde{\mu}_t(x_t) &= \frac{1}{\sqrt{\alpha_t}} \left(x_t - \frac{1 - \alpha_t}{\sqrt{1 - \bar{\alpha}_t}} \varepsilon_t \right)
\end{aligned}$$

$$\begin{aligned}
p_\theta(x_{t-1} | x_t) &:= \mathcal{N}(x_{t-1}; \boldsymbol{\mu}_\theta(x_t, t), \Sigma_\theta(x_t, t)) \quad \text{Predicted mean with neural network} \\
p_\theta(x_{t-1} | x_t) &:= \mathcal{N}(x_{t-1}; \boldsymbol{\mu}_\theta(x_t, t), \beta I) \quad \text{Beta is fixed}
\end{aligned}$$

Initially, the neural network approximated the mean.

$$\begin{aligned}
\tilde{\boldsymbol{\mu}}_t(x_t) &= \frac{1}{\sqrt{\alpha_t}} \left(x_t - \frac{1 - \alpha_t}{\sqrt{1 - \bar{\alpha}_t}} \varepsilon_t \right) \\
\boldsymbol{\mu}_\theta(x_t, t) &= \frac{1}{\sqrt{\alpha_t}} \left(x_t - \frac{1 - \alpha_t}{\sqrt{1 - \bar{\alpha}_t}} \varepsilon_\theta(x_t, t) \right)
\end{aligned}$$

The authors use a simple Mean Squared Error to compare the actual mean and predicted mean.

$$\begin{aligned}
L_t &= \mathbb{E}_{x_0, \varepsilon} \left[\frac{1}{2\sigma_t^2} \|\tilde{\boldsymbol{\mu}}_t(x_t) - \boldsymbol{\mu}_\theta(x_t, t)\|^2 \right] \\
&= \mathbb{E}_{x_0, \varepsilon} \left[\frac{1}{2\sigma_t^2} \left\| \frac{1}{\sqrt{\alpha_t}} \left(x_t - \frac{1 - \alpha_t}{\sqrt{1 - \bar{\alpha}_t}} \varepsilon_t \right) - \frac{1}{\sqrt{\alpha_t}} \left(x_t - \frac{1 - \alpha_t}{\sqrt{1 - \bar{\alpha}_t}} \varepsilon_\theta(x_t, t) \right) \right\|^2 \right] \\
&= \mathbb{E}_{x_0, \varepsilon} \left[\frac{(1 - \alpha_t)^2}{2\alpha_t(1 - \bar{\alpha}_t)\sigma_t^2} \|\varepsilon_t - \varepsilon_\theta(x_t, t)\|^2 \right] \\
L_{\text{simple}}(\theta) &:= \mathbb{E}_{t, \mathbf{x}_0, \epsilon} \left[\left\| \epsilon - \epsilon_\theta(\sqrt{\alpha_t} \mathbf{x}_0 + \sqrt{1 - \bar{\alpha}_t} \epsilon, t) \right\|^2 \right] \quad \text{Disregard weighted term}
\end{aligned}$$

Authors use simple MSE without the weighted term because the quality of samples is better. It is unnecessary to use the weighted equation. Standard L2 loss with predicted noise and actual noise works better.

L_0 - Reconstruction Term - This term is thrown away by the authors. This term can be approximated by the same neural network in the previous term. Ignoring this term makes sampling better.

$$\begin{aligned}
p_\theta(\mathbf{x}_0 \mid \mathbf{x}_1) &= \prod_{i=1}^D \int_{\delta_-(x_0^i)}^{\delta_+(x_0^i)} \mathcal{N}(x; \mu_\theta^i(\mathbf{x}_1, 1), \sigma_1^2) dx \\
\delta_+(x) &= \begin{cases} \infty & \text{if } x = 1 \\ x + \frac{1}{255} & \text{if } x < 1 \end{cases} \quad \delta_-(x) = \begin{cases} -\infty & \text{if } x = -1 \\ x - \frac{1}{255} & \text{if } x > -1 \end{cases}
\end{aligned} \tag{13}$$

3.2 GitHub Repository Link

<https://github.com/CamdenKitowski/DenoisingDiffusionModels>

4 Math Terms and Definitions

Variational Inference - the method/approach to approximate complex probability distribution P with a simpler one Q . This is used when exact inference is intractable. Usually done by minimizing KL divergence between Q and P .

Evidence Lower Bound (ELBO) - objective function (function that is either maximized or minimized to solve a problem) used in variational inference. It represents a lower bound on the log-likelihood of the observed data.

Two part objective:

- Maximize expected log-likelihood of data under the approximate distribution Q
- Minimize KL divergence between Q and prior distribution over the latent variables

Kullback-Leibler (KL) divergence - known as relative entropy - measure of how one probability distribution diverges from another one.

For discrete variables:

$$D_{KL}(P \parallel Q) = \sum_x P(x) \log \left(\frac{P(x)}{Q(x)} \right)$$

For continuous variables

$$D_{KL}(P||Q) = \int_{-\infty}^{\infty} p(x) \log \left(\frac{p(x)}{q(x)} \right) dx$$

where $p(x)$ and $q(x)$ are probability density functions of P and Q , respectively. The divergence is always non-negative. If it is zero, then P and Q are the same distribution.

Likelihood and Log-Likelihood -

can you explain why we start with negative log likelihood?

Likelihood - the likelihood of a set of parameters θ given observations X is the probability of X given θ , denoted $P(X|\theta)$. The likelihood function $L(\theta|X)$ is defined as the joint probability/density of the observed data as a function of the parameters θ . Use this to find the parameters values that maximize the probability. Known as Maximum Likelihood Estimation

Log-Likelihood - Use the natural logarithm of likelihood function, denoted as $\log(L(\theta|X))$. Turns products into sums, making it easier to differentiate and more stable when looking for the maximum likelihood estimates.

More Notes - Both serve the same purpose. Maximizing the likelihood is equivalent to maximizing the log-likelihood because log is strictly increasing function. So the set of parameters that maximizes the likelihood also maximizes the log-likelihood.

Markov Chain - model systems that transition from one state to another, where the probability of each subsequent state depends only on the current state and not on the sequence of events that preceded it

Expected Value - measure of central tendency of its probability distribution. Calculated by summing all possible values of the variable, each multiplied by its probability of occurrence

$$E[X] = \sum_{i=1}^n x_i P(X = x_i)$$

Conditional Mean/Expectation - it's the expected value of a random variable given that a certain condition(s) holds true.

$$E[Y|X = x] = \sum_y y \cdot P(Y = y|X = x)$$

Probability vs. Expected Value -

- $P(y|x)$ is used when we're interested in the chance of an event happening.

- $E(y|x)$ is used in situations where we want to predict or estimate the average outcome of a random variable given some condition or information.

Forward Process Equation vs. Forward Process Posteriors -

Forward Process Equation - tells you how to go from one time step to the next by adding noise. It takes the form of a Markov Chain in the form of Gaussian transitions.

Forward Process Posteriors - Describe the probability distribution over possible states at each time step that is conditioned on the original data. Expressed in terms of the posteriors parameters which is mean and variance.

More Notes - The equation provides a procedural generation path, and the posteriors provide a statistical description of where you might end up along that path.