



250830

AAILAB
Department of Industrial and Systems Engineering
KAIST

Problem setup

- Adjacency matrix
 - A^* : adjacency matrix for default case.
 - $A'_{(1)}, \dots, A'_{(K)}$: adjacency matrix for K exceptional case.
- Path dataset
 - We can access i.i.d. samples from $p_{A^*}(x)$.
 - We can access i.i.d. samples from $p_{A'_{(k)}}(x)$ for $k \in \{1, \dots, K\}$.
 - A path $x = \{v_1, \dots, v_{|x|}\}$ consists of $|x|$ number of vertices.
- Goal
 - Generating samples from $p_{A'_{\text{test}}}(x)$ for a given adjacency matrix A'_{test} .

- Motivation
 - In the default case (A^*), we usually have access to larger number of data, so a base diffusion model is trained with a focus on generating high-quality data.
 - We train a case-agnostic auxiliary network to correct the distribution ($A^* \rightarrow A'_k$) based on limited data from various exceptions, enabling adaptation to unseen cases at test time.
- 1. Training a base diffusion model that approximates $p_{A^*}(x)$ corresponds to a default case.
 - Forward process
 - $q_{A^*}(x_t|x_0) = \bigotimes_{i=1}^{|x|} Cat(x_t^i; p = x_0^i \bar{C}_{t-1}^*)$, where $C_t^* = \exp\{(A^* - D^*)t\}$, $\bar{C}_t^* = C_1^* \cdots C_t^*$.
 - Reverse process
 - $q_{A^*}(x_{t-1}|x_t) = \bigotimes_{i=1}^{|x|} Cat(x_{t-1}^i; p \propto x_t^i C_t^* \odot \hat{x}_0^i \bar{C}_{t-1}^*)$, where $\hat{x}_0 = NN(x_t, t)$.
- 2. Training a discriminator that approximates $\frac{p_{A'}(x)}{p_{A^*}(x)}$.
 - Learn a distribution correction strategy within the denoising diffusion process.
 - We need to estimating $\frac{q_{A'}(x_{t-1}|x_t)}{q_{A^*}(x_{t-1}|x_t)}$.

Discriminator design

- Denoising probability with base diffusion process.
 - $q_{A^*}(x_{t-1}|x_t) = \bigotimes_{i=1}^{|x|} \text{Cat}(x_{t-1}^i; p \propto x_t^i C_t^* \odot \hat{x}_0^i \bar{C}_{t-1}^*)$ where $C_t^* = \exp\{(A^* - D^*)t\}$, $\bar{C}_t^* = C_1^* \cdots C_t^*$
 → This calculate probability for all possible states of $x_{t-1}^i \in V$ and for all i with single feed-forward.
- Target correction term that needs to estimate.
 - $\frac{q_{A'}(x_{t-1}|x_t)}{q_{A^*}(x_{t-1}|x_t)} \approx \bigotimes_{i=1}^{|x|} \frac{q_{A'}(x_{t-1}^i|x_t)}{q_{A^*}(x_{t-1}^i|x_t)} = \boxed{\bigotimes_{i=1}^{|x|} \frac{P(A = A'|x_{t-1}^i, x_t)}{P(A = A^*|x_{t-1}^i, x_t)}}$ for all possible state of $x_{t-1}^i \in V$ and for all i .
- Discriminator function: discriminates the samples from $q_{A^*}(x_t)$ and the samples from $q_{A'}(x_t)$.
 - $D_\phi: (A' \times x_t \times t) \rightarrow P_\phi(A = A'|x_t, A \in \{A^*, A'\}, t)$
 - The target adjacency matrix A' becomes an input to the discriminator to consider arbitrary testing situations.
 - $\frac{P(A = A'|x_t)}{P(A = A^*|x_t)} \approx \frac{D_\phi(x_t, A', t)}{1 - D_\phi(x_t, A', t)} = \boxed{\frac{P_\phi(A = A'|x_t)}{P_\phi(A = A^*|x_t)}}$ → Challenge: how can the discriminator represent target correction term with single feed-forward?

Approximation for tractable discriminator guidance

Target: $\otimes_{i=1}^{|x|} \frac{P(A = A' | x_{t-1}^i, x_t)}{P(A = A^* | x_{t-1}^i, x_t)}$ for all possible state of $x_{t-1}^i \in V$ and for all i .

We can compute: $\frac{P_\phi(A = A' | x_t)}{P_\phi(A = A^* | x_t)}$

- Approximation 1. the target denoising vertex as the vertex at the current time.
 - $P(A = A' | x_{t-1}^i, x_t) \approx P(A = A' | \tilde{x}_t^i)$, where $\tilde{x}_t^i = [x_t^{(1:i-1)}, x_{t-1}^i, x_t^{(i+1:|x|)}]$.
 - This make sure that only the path from single time step is given to the discriminator.
- Approximation 2. first order approximation for single feed-forward inference
 - $P(A = A' | \tilde{x}_t^i) = \exp((\tilde{x}_t^i - x_t) \nabla \log P(A = A' | x_t) + \log P(A = A' | x_t))$
 - The neural estimation target $P(A = A' | x_t)$ is computed independently of the denoising state $x_{t-1}^i \in V$ or any specific vertex i . This means we can get the categorical probabilities for all the vertices needed for denoising with just one discriminator inference and simple tensor operations.

This approximation applies the method from (Schiff et al., 2025).

(Schiff et al., 2025) Simple Guidance Mechanisms for Discrete Diffusion Models, ICLR 2025.

Discriminator training algorithm

- Training algorithm for A' conditioned discriminator. (consider $k = 2$)
 - We have discriminator $(A' \times x_t \times t) \rightarrow P_\phi(A = A'|x_t, A \in \{A^*, A'\}, t)$
 - We can access i.i.d. samples from $\mathcal{D}_{A^*} = \{x_{00}^{*i}\}_{i=1}^N$, $\mathcal{D}_{A'_1} = \{x_{10}'^i\}_{i=1}^N$, $\mathcal{D}_{A'_2} = \{x_{20}'^i\}_{i=1}^N$
1. Get $\{x_{00}^{*i}\}_{i=1}^B$, $\{x_{10}'^i\}_{i=1}^{B/2}$, and $\{x_{20}'^i\}_{i=1}^{B/2}$ and perturb (A^*, A'_1, A'') them into $\{x_{0t}^{*i}, 0\}_{i=1}^B$, $\{x_{1t}'^i, 1\}_{i=1}^{B/2}$, $\{x_{2t}'^i, 1\}_{i=1}^{B/2}$.
 2. For $\{x_{0t}^{*i}\}_{i=1}^{B/2}$, $\{x_{1t}'^i\}_{i=1}^{B/2}$, compute $D_\phi(A'_1, v_t, t)$ and compute BCE loss with its label {0,1}.
 3. For $\{x_{0t}^{*i}\}_{i=B/2}^B$, $\{x_{2t}'^i\}_{i=1}^{B/2}$, compute $D_\phi(A'_2, v_t, t)$ and compute BCE loss with its label {0,1}.