

Diffusion Models for Time Series Applications: A Survey

Lequan Lin, Zhengkun Li, Ruikun Li, Xuliang Li, and Junbin Gao

Discipline of Business Analytics, The University of Sydney Business School

The University of Sydney, Camperdown, NSW 2006, Australia

{lequan.lin, zhengkun.li, ruikun.li, xuliang.li, junbin.gao}@sydney.edu.au

Abstract

Diffusion models, a family of generative models based on deep learning, have become increasingly prominent in cutting-edge machine learning research. With a distinguished performance in generating samples that resemble the observed data, diffusion models are widely used in image, video, and text synthesis nowadays. In recent years, the concept of diffusion has been extended to time series applications, and many powerful models have been developed. Considering the deficiency of a methodical summary and discourse on these models, we provide this survey as an elementary resource for new researchers in this area and also an inspiration to motivate future research. For better understanding, we include an introduction about the basics of diffusion models. Except for this, we primarily focus on diffusion-based methods for time series forecasting, imputation, and generation, and present them respectively in three individual sections. We also compare different methods for the same application and highlight their connections if applicable. Lastly, we conclude the common limitation of diffusion-based methods and highlight potential future research directions.

1 Introduction

Diffusion models, a family of deep learning-based generative models, have risen to prominence in the machine learning community in recent years (Croitoru et al., 2023; Yang et al., 2022a). With exceptional performance in various real-world applications such as image synthesis (Austin et al., 2021; Dhariwal and Nichol, 2021; Ho et al., 2022a), video generation (Harvey et al., 2022; Ho et al., 2022b; Yang et al., 2022b), natural language processing (Li et al., 2022; Nikolay et al., 2022; Yu et al., 2022), and time series prediction (Rasul et al., 2021a; Li et al., 2022; Alcaraz and Strodthoff, 2023), diffusion models have demonstrated their power over many existing generative techniques.

Given some observed data \mathbf{x} from a target distribution $q(\mathbf{x})$, the objective of a generative model is to learn a generative process that produces new samples from $q(\mathbf{x})$

(Luo, 2022). To learn such a generative process, most diffusion models begin with progressively disturbing the observed data by injecting Gaussian noises, then applying a reversed process with a learnable transition kernel to recover the data (Sohl-Dickstein et al., 2015; Ho et al., 2020; Luo, 2022). Typical diffusion models assume that after a certain number of noise injection steps, the observed data will become standard Gaussian noises. So, if we can find the probabilistic process that recovers the original data from standard Gaussian noises, then we can generate similar samples using the same probabilistic process with any random standard Gaussian noises as the starting point.

The recent three years have witnessed the extension of diffusion models to time series-related applications, including time series forecasting (Rasul et al., 2021a; Li et al., 2022; Biloš et al., 2022), time series imputation (Tashiro et al., 2021; Alcaraz and Strodthoff, 2023; Liu et al., 2023), and time series generation (Lim et al., 2023). Given observed historical time series, we often try to predict future time series. This process is known as time series forecasting. Since observed time series are sometimes incomplete due to reasons such as data collection failures and human errors, time series imputation is implemented to fill in the missing values. Different from time series forecasting and imputation, time series generation or synthesis aims to produce more time series samples with similar characteristics as the observed period.

Basically, diffusion-based methods for time series applications are developed from three fundamental formulations, including denoising diffusion probabilistic models (DDPMs), score-based generative models (SGMs), and stochastic differential equations (SDEs). The target distributions learned by the diffusion components in different methods often involve the condition on previous time steps. Nevertheless, the design of the diffusion and denoising processes varies with different objectives of different tasks. Hence, a comprehensive and self-contained summary of relevant literature will be an inspiring beacon for new researchers who just enter this new-born area and experienced researchers who seek for future directions. Accordingly, this survey aims to summarize existing literature, compare different approaches, and identify potential limitations.

In this paper, we will review on diffusion-based models for time series applications (please refer to **Table 1** for a quick summary). For a better understanding, we will include a brief introduction about three predominant formulations of diffusion models in section 2. Next, we will categorize the existing models based on their major functions. More specifically, we will discuss the models primarily for time series forecasting, time series imputation, and time series generation in section 3, section 4, and section 5, respectively. In each section, we will have a separate subsection for problem formulation, which helps to clarify the objective, training and forecasting settings of each specific application. We will highlight if a model can serve multiple purposes and articulate the linkage when one model is related to or slightly different from another. Eventually, we will conclude this survey in section 6.

2 Basics of Diffusion Models

The underlying principle of diffusion models is to progressively perturb the observed data with a forward diffusion process, then recover the original data through a backward

Table 1: A summary of diffusion-based methods for time series applications.

Application	Data Type	Method	Source
Time Series Forecasting	Multivariate Time Series	TimeGrad	Rasul et al. (2021a)
		ScoreGrad	Yan et al. (2021)
		D ³ VAE	Li et al. (2022)
		DSPD/CSPD	Biloš et al. (2022)
	Spatio-temporal Graphs	DiffSTG	Wen et al. (2023)
		GCRDD	Li et al. (2023)
Time Series Imputation	Multivariate Time Series	CSDI	Tashiro et al. (2021)
		DSPD/CSPD	Biloš et al. (2022)
		SSSD	Alcaraz and Strodthoff (2023)
	Spatio-temporal Graphs	PriSTI	Liu et al. (2023)
Time Series Generation	Multivariate Time Series	TSGM	Lim et al. (2023)

reverse process. The forward process involves multiple steps of noise injection, where the noise level changes at each step. The backward process, on the contrary, consists of multiple denoising steps that aim to remove the injected noises gradually. Normally, the backward process is parameterized by a neural network. Once the backward process has been learned, it can generate new samples from almost arbitrary initial data. Stemming from this basic idea, diffusion models are predominantly formulated in three ways: denoising diffusion probabilistic models (DDPMs), score-based generative models (SGMs), and stochastic differential equations (SDEs). In this section, we will briefly review on these three formulations of diffusion models.

2.1 Denoising Diffusion Probabilistic Models

DDPMs implement the forward and backward processes through two Markov chains (Sohl-Dickstein et al., 2015; Ho et al., 2020). Let the original observed data be \mathbf{x}^0 , where 0 indicates that the data are free from the noises injected in the diffusion process.

The forward Markov chain transforms \mathbf{x}^0 to a sequence of disturbed data $\mathbf{x}^1, \mathbf{x}^2, \dots, \mathbf{x}^K$ with a diffusion transition kernel:

$$q(\mathbf{x}^k | \mathbf{x}^{k-1}) = \mathcal{N}\left(\mathbf{x}^k; \sqrt{\alpha_k} \mathbf{x}^{k-1}, (1 - \alpha_k) \mathbf{I}\right), \quad (1)$$

where $\alpha_k \in (0, 1)$ for $k = 1, 2, \dots, K$ are hyperparameters indicating the changing variance of the noise level at each step, and $\mathcal{N}(\mathbf{x}; \boldsymbol{\mu}, \boldsymbol{\Sigma})$ is the general notation for the Gaussian distribution of \mathbf{x} with the mean $\boldsymbol{\mu}$ and the covariance $\boldsymbol{\Sigma}$, respectively. A nice property of this Gaussian transition kernel is that we may obtain \mathbf{x}^k directly from \mathbf{x}^0 by

$$q(\mathbf{x}^k | \mathbf{x}^0) = \mathcal{N}\left(\mathbf{x}^k; \sqrt{\tilde{\alpha}_k} \mathbf{x}^0, (1 - \tilde{\alpha}_k) \mathbf{I}\right), \quad (2)$$

where $\tilde{\alpha}_k := \prod_{i=1}^k \alpha_i$. Therefore, $\mathbf{x}^k = \sqrt{\tilde{\alpha}_k} \mathbf{x}^0 + \sqrt{1 - \tilde{\alpha}_k} \boldsymbol{\epsilon}$ with $\boldsymbol{\epsilon} \sim \mathcal{N}(\mathbf{0}, \mathbf{I})$. Normally, we design $\tilde{\alpha}_K \approx 0$ such that $q(\mathbf{x}^K) := \int q(\mathbf{x}^K | \mathbf{x}^0) q(\mathbf{x}^0) d\mathbf{x}^0 \approx \mathcal{N}(\mathbf{x}^K; \mathbf{0}, \mathbf{I})$, which means the starting point of the backward chain can be any standard Gaussian noises.

The reverse transition kernel is modelled by a parameterized neural network

$$p_{\boldsymbol{\theta}}(\mathbf{x}^{k-1} | \mathbf{x}^k) = \mathcal{N}\left(\mathbf{x}^{k-1}; \boldsymbol{\mu}_{\boldsymbol{\theta}}(\mathbf{x}^k, k), \boldsymbol{\Sigma}_{\boldsymbol{\theta}}(\mathbf{x}^k, k)\right), \quad (3)$$

where $\boldsymbol{\theta}$ denotes learnable parameters. Now, the remaining problem is how to estimate $\boldsymbol{\theta}$. Basically, the objective is to maximize the likelihood objective function so that the probability of observing the training sample \mathbf{x}^0 estimated by $p_{\boldsymbol{\theta}}(\mathbf{x}^0)$ is maximized. This task is accomplished by minimizing the variational lower bound of the estimated negative log-likelihood $\mathbb{E}[-\log p_{\boldsymbol{\theta}}(\mathbf{x}^0)]$, that is,

$$\mathbb{E}_{q(\mathbf{x}^{0:K})} \left[-\log p(\mathbf{x}^K) - \sum_{k=1}^K \log \frac{p_{\boldsymbol{\theta}}(\mathbf{x}^{k-1}|\mathbf{x}^k)}{q(\mathbf{x}^k|\mathbf{x}^{k-1})} \right] \quad (4)$$

where $\mathbf{x}^{0:K}$ denotes the sequence $\mathbf{x}^0, \dots, \mathbf{x}^K$.

Ho et al. (2020) proposed that we could simplify the covariance matrix $\boldsymbol{\Sigma}_{\boldsymbol{\theta}}(\mathbf{x}^k, k)$ in Equation (3) as a constant-dependent matrix $\sigma_k^2 \mathbf{I}$, where σ_k^2 controls the noise level and may vary at different diffusion steps. Besides, they rewrote the mean as a function of a learnable noise term as:

$$\boldsymbol{\mu}_{\boldsymbol{\theta}}(\mathbf{x}^k, k) = \frac{1}{\sqrt{\alpha_k}} \left(\mathbf{x}^k - \zeta(k) \boldsymbol{\epsilon}_{\boldsymbol{\theta}}(\mathbf{x}^k, k) \right), \quad (5)$$

where $\zeta(k) = \frac{1-\alpha_k}{\sqrt{1-\bar{\alpha}_k}}$, and $\boldsymbol{\epsilon}_{\boldsymbol{\theta}}$ is a noise-matching network that predicts $\boldsymbol{\epsilon}$ corresponding to inputs \mathbf{x}^k and k . With the property in Equation (2), Ho et al. (2020) further simplifies the objective function to

$$\mathbb{E}_{k, \mathbf{x}^0, \boldsymbol{\epsilon}} \left[\delta(k) \left\| \boldsymbol{\epsilon} - \boldsymbol{\epsilon}_{\boldsymbol{\theta}} \left(\sqrt{\bar{\alpha}_k} \mathbf{x}^0 + \sqrt{1-\bar{\alpha}_k} \boldsymbol{\epsilon}, k \right) \right\|^2 \right], \quad (6)$$

where $\delta(k) = \frac{(1-\alpha_k)^2}{2\sigma_k^2 \alpha_k (1-\bar{\alpha}_k)}$ is a positive-valued weight that can be discarded to produce better performance in practice.

Eventually, samples are generated by eliminating the noises in $\mathbf{x}^K \sim \mathcal{N}(\mathbf{x}^K; \mathbf{0}, \mathbf{I})$. More specifically, for $k = K-1, K-2, \dots, 0$,

$$\mathbf{x}^k \leftarrow \frac{(\mathbf{x}^{k+1} - \zeta(k+1) \boldsymbol{\epsilon}_{\boldsymbol{\theta}}(\mathbf{x}^{k+1}, k+1))}{\sqrt{\alpha_{k+1}}} + \sigma_k \mathbf{z},$$

where $\mathbf{z} \sim \mathcal{N}(\mathbf{0}, \mathbf{I})$ for $k = K-1, \dots, 1$, and $\mathbf{z} = \mathbf{0}$ for $k = 0$.

2.2 Score-based Generative Models

Score-based generative models (SGMs) consist of two modules, including score matching and annealed Langevin dynamics (ALD). ALD is a sampling algorithm that generates samples with an iterative process by applying Langevin Monte Carlo at each update step (Song and Ermon, 2019). Stein score is an essential component of ALD. The Stein score of a density function $q(\mathbf{x})$ is defined as $\nabla_{\mathbf{x}} \log q(\mathbf{x})$. Since the true probabilistic distribution $q(\mathbf{x})$ is usually unknown, score matching (Hyvärinen and Dayan, 2005) is implemented to approximate the Stein score with a score-matching network. Here we primarily focus on denoising score matching (Vincent, 2011) as it is empirically more efficient, but other methods such as sliced score matching (Song et al., 2020) are also commonly mentioned in the relevant literature.

The underlying principle of denoising score matching is to process the observed data with the forward transition kernel $q(\mathbf{x}^k|\mathbf{x}^0) = \mathcal{N}(\mathbf{x}^k; \mathbf{x}^0, \sigma_k^2 \mathbf{I})$, with σ_k^2 being a set of increasing noise levels for $k = 1, \dots, K$, and then jointly estimate the Stein scores for the noise density distributions $q_{\sigma_1}(\mathbf{x}), q_{\sigma_2}(\mathbf{x}), \dots, q_{\sigma_K}(\mathbf{x})$ (Song and Ermon, 2019). The Stein score for noise density function $q_{\sigma_k}(\mathbf{x})$ is defined as $\nabla_{\mathbf{x}} \log q_{\sigma_k}(\mathbf{x})$.

Then, the Stein score is approximated by a neural network $\mathbf{s}_{\theta}(\mathbf{x}, \sigma_k)$, where θ contains learnable parameters. Accordingly, the initial objective function is given as

$$\mathbb{E}_{k, \mathbf{x}^0, \mathbf{x}^k} \left[\left\| \mathbf{s}_{\theta}(\mathbf{x}^k, k) - \nabla_{\mathbf{x}^k} \log q_{\sigma_k}(\mathbf{x}^k) \right\|^2 \right]. \quad (7)$$

With the Gaussian assumption of the forward transition kernel, a tractable version of the objective function can be found as

$$\mathbb{E}_{k, \mathbf{x}^0, \mathbf{x}^k} \left[\delta(k) \left\| \mathbf{s}_{\theta}(\mathbf{x}^k, \sigma_k) + \frac{\mathbf{x}^k - \mathbf{x}^0}{\sigma_k^2} \right\|^2 \right], \quad (8)$$

where $\delta(k)$ is a positive-valued weight depending on the noise scale σ_k .

After the score-matching network \mathbf{s}_{θ} is learned, the ALD algorithm will be implemented for sampling. The algorithm is initialized with a sequence of increasing noise levels $\sigma_1, \dots, \sigma_K$ and a starting point $\mathbf{x}^{K,0} \sim \mathcal{N}(\mathbf{0}, \mathbf{I})$. For $k = K, K-1, \dots, 0$, \mathbf{x}^k will be updated with N iterations that compute

$$\begin{aligned} \mathbf{z} &\leftarrow \mathcal{N}(\mathbf{0}, \mathbf{I}) \\ \mathbf{x}^{k,n} &\leftarrow \mathbf{x}^{k,n-1} + \frac{1}{2} \eta_k \mathbf{s}_{\theta}(\mathbf{x}^{k,n-1}, \sigma_k) + \sqrt{\eta_k} \mathbf{z}, \end{aligned}$$

where $n = 1, \dots, N$, $\mathbf{z} \sim \mathcal{N}(\mathbf{0}, \mathbf{I})$, and η_k represents the step of update. Note that after each N iterations, the last output $\mathbf{x}^{k,N}$ will be assigned as the starting point of the next N iterations, that is, $\mathbf{x}^{k-1,1}$. $\mathbf{x}^{0,N}$ will be the final sample. The role of \mathbf{z} in this sampling process is to slightly add uncertainty such that the algorithm will not end up with almost identical samples.

2.3 Stochastic Differential Equations

DDPMs and SGMs implement the forward pass as a discrete process, which means we should carefully design the diffusion steps. To overcome this limitation, one may consider the diffusion process as continuous such that it becomes the solution of a stochastic differential equation (SDE) (Song et al., 2021). This formulation can be thought of as a generalization of the previous two formulations since both DDPMs and SGMs are discrete forms of SDEs. The backward process is modelled as a time-reverse SDE, and samples can be generated by solving this time-reverse SDE. Let \mathbf{w} and $\tilde{\mathbf{w}}$ be a standard Wiener process and its time-reverse version, respectively, and consider a continuous diffusion time $k \in [0, K]$. A general expression of SDE is

$$d\mathbf{x} = f(\mathbf{x}, k)dk + g(k)d\mathbf{w}, \quad (9)$$

and the time-reverse SDE, as shown by Anderson (1982), is

$$d\mathbf{x} = [f(\mathbf{x}, k) - g(k)^2 \nabla_{\mathbf{x}} \log q_k(\mathbf{x})] dk + g(k)d\tilde{\mathbf{w}}, \quad (10)$$

In addition, Song et al. (2021) has illustrated that sampling from the probability flow ordinary differential equation (ODE) as following has the same distribution as the time-reverse SDE:

$$d\mathbf{x} = \left[f(\mathbf{x}, k) - \frac{1}{2}g(k)^2 \nabla_{\mathbf{x}} \log q_k(\mathbf{x}) \right] dk. \quad (11)$$

Here $f(\mathbf{x}, k)$ and $g(k)$ separately compute the drift coefficient and the diffusion coefficient for the diffusion process. $\nabla_{\mathbf{x}} \log q_k(\mathbf{x})$ is the Stein score corresponding to the marginal distribution of \mathbf{x}^k , which is unknown but can be learned with a similar method as in SGMs with the objective function

$$\mathbb{E}_{k, \mathbf{x}^0, \mathbf{x}^k} \left[\delta(k) \left\| \mathbf{s}_{\theta}(\mathbf{x}^k, k) - \nabla_{\mathbf{x}^k} \log q_{0k}(\mathbf{x}^k | \mathbf{x}^0) \right\|^2 \right]. \quad (12)$$

Now, how to write the diffusion processes of DDPMs and SGMs as SDEs? Recall that α_k is a defined parameter in DDPMs and σ_k^2 denotes the noise level in SGMs. The SDE corresponding to DDPMs is known as variance preserving (VP) SDE, defined as

$$d\mathbf{x} = -\frac{1}{2}\alpha(k)\mathbf{x}dk + \sqrt{\alpha(k)}d\mathbf{w}, \quad (13)$$

where $\alpha(\cdot)$ is a continuous function, and $\alpha\left(\frac{k}{K}\right) = K(1 - \alpha_k)$ as $K \rightarrow \infty$. For the forward pass of SGMs, the associated SDE is known as variance exploding (VE) SDE, defined as

$$d\mathbf{x} = \sqrt{\frac{d[\sigma(k)^2]}{dk}}d\mathbf{w}, \quad (14)$$

where $\sigma(\cdot)$ is a continuous function, and $\sigma\left(\frac{k}{K}\right) = \sigma_k$ as $K \rightarrow \infty$ (Song et al., 2021). Inspired by VP SDE, Song et al. (2021) designed another SDE called sub-VP SDE that performs especially well on likelihoods, given by

$$d\mathbf{x} = -\frac{1}{2}\alpha(k)\mathbf{x}dk + \sqrt{\alpha(k) \left(1 - e^{-2\int_0^k \alpha(s)ds}\right)}d\mathbf{w}. \quad (15)$$

The objective function involves a perturbation distribution $q_{0k}(\mathbf{x}^k | \mathbf{x}^0)$ that varies for different SDEs. For the three aforementioned SDEs, their corresponding perturbation distributions are derived as

$$q_{0k}(\mathbf{x}^k | \mathbf{x}^0) = \begin{cases} \mathcal{N}(\mathbf{x}^k; \mathbf{x}^0, [\sigma(k)^2 - \sigma(0)^2]\mathbf{I}), & \text{(VP SDE)} \\ \mathcal{N}(\mathbf{x}^k; \mathbf{x}^0 e^{-\frac{1}{2}\int_0^k \alpha(s)ds}, [1 - e^{-\int_0^k \alpha(s)ds}]\mathbf{I}), & \text{(VE SDE)} \\ \mathcal{N}(\mathbf{x}^k; \mathbf{x}^0 e^{-\frac{1}{2}\int_0^k \alpha(s)ds}, [1 - e^{-\int_0^k \alpha(s)ds}]^2\mathbf{I}), & \text{(sub-VP SDE)} \end{cases} \quad (16)$$

After successfully learning $\mathbf{s}_{\theta}(\mathbf{x}, k)$, samples are produced by deriving the solutions to the time-reverse SDE or the probability flow ODE with techniques such as ALD.

3 Time Series Forecasting

Multivariate time series forecasting is a crucial area of study in machine learning research, with wide-ranging applications across a variety of industries. Different from univariate

time series, which only track one feature over time, multivariate time series involve the historical observations of multiple features that interact with each other and evolve with time. Consequently, they provide a more comprehensive understanding of complex systems and realize more reliable predictions of future trends and behaviours.

In recent years, generative models have been implemented for multivariate time series forecasting tasks. For example, WaveNet is a generative model with dilated causal convolutions that encode long-term dependencies for sequence prediction (Oord et al., 2016). As another example, Rasul et al. (2021b) model multivariate time series with an autoregressive deep learning model, in which the data distribution is expressed by a conditional normalizing flow. Nevertheless, the common shortcoming of these models is that the functional structure of their target distributions are strictly constrained. Diffusion-based methods, on the other hand, can provide a less restrictive solution. In this section, we will discuss four diffusion-based approaches. We also include two models designed specifically for spatio-temporal graphs (i.e., spatially related entities with multivariate time series) to highlight the extension of diffusion theories to more complicated problem settings. Since relevant literature mostly focuses on multivariate time series forecasting, “forecasting” refers to multivariate time series forecasting in the rest of this survey unless otherwise stated.

3.1 Problem Formulation

Consider a multivariate time series $\mathbf{X}^0 = \{\mathbf{x}_1^0, \mathbf{x}_2^0, \dots, \mathbf{x}_T^0 | \mathbf{x}_i^0 \in \mathbb{R}^D\}$, where 0 indicates that the data is free from the perturbation in the diffusion process. The forecasting task is to predict $\mathbf{X}_p^0 = \{\mathbf{x}_{t_0}^0, \mathbf{x}_{t_0+1}^0, \dots, \mathbf{x}_T^0\}$ given the historical information $\mathbf{X}_c^0 = \{\mathbf{x}_1^0, \mathbf{x}_2^0, \dots, \mathbf{x}_{t_0-1}^0\}$. \mathbf{X}_c^0 is known as the context window, while \mathbf{X}_p^0 is known as the prediction interval. In diffusion-based models, the problem is formulated as learning the joint probabilistic distribution of data in the prediction interval:

$$q(\mathbf{x}_{t_0:T}^0 | \mathbf{x}_{1:t_0-1}^0) = \prod_{t=t_0}^T q(\mathbf{x}_t^0 | \mathbf{x}_{1:t_0-1}^0). \quad (17)$$

Some literature also considers the role of covariates in forecasting, such as (Rasul et al., 2021a) and (Yan et al., 2021). Covariates are additional information that may impact the behaviour of variables over time, such as seasonal fluctuations and weather changes. Incorporating covariates in forecasting often helps to strengthen the identification of factors that drive temporal trends and patterns in data. The forecasting problem with covariates is formulated as

$$q(\mathbf{x}_{t_0:T}^0 | \mathbf{x}_{1:t_0-1}^0, \mathbf{c}_{1:T}) = \prod_{t=t_0}^T q(\mathbf{x}_t^0 | \mathbf{x}_{1:t_0-1}^0, \mathbf{c}_{1:T}), \quad (18)$$

where $\mathbf{c}_{1:T}$ denotes the covariates for all time points and is assumed to be known for the whole period.

For the purpose of training, one may randomly sample the context window followed by the prediction window from the complete training data. This process can be seen as applying a moving window with size T on the whole timeline. Then, the optimization of

the objective function can be conducted with the samples. Forecasting future time series is usually achieved by the generation process corresponding to the diffusion models.

3.2 TimeGrad

The first noticeable work on diffusion-based forecasting is TimeGrad proposed by Rasul et al. (2021a). Developed from DDPM models, TimeGrad firstly injects noises to data at each predictive time point, and then gradually denoise through a backward transition kernel conditioned on historical time series. To encode historical information, TimeGrad approximate the conditional distribution in Equation (18) by

$$\prod_{t=t_0}^T p_{\theta}(\mathbf{x}_t^0 | \mathbf{h}_{t-1}), \quad (19)$$

where

$$\mathbf{h}_t = \text{RNN}_{\theta}(\mathbf{x}_t^0, \mathbf{c}_t, \mathbf{h}_{t-1}) \quad (20)$$

is the hidden state calculated with a RNN module such as LSTM (Hochreiter and Schmidhuber, 1997) or GRU (Chung et al., 2014) that can preserve historical temporal information, and θ contains learnable parameters for the overall conditional distribution and its RNN component.

The objective function of TimeGrad is in the form of a negative log-likelihood, given as

$$\sum_{t=t_0}^T -\log p_{\theta}(\mathbf{x}_t^0 | \mathbf{h}_{t-1}), \quad (21)$$

where for each $t \in [t_0, T]$, $-\log p_{\theta}(\mathbf{x}_t^0 | \mathbf{h}_{t-1})$ is upper bounded by

$$\mathbb{E}_{k, \mathbf{x}_t^0, \epsilon} \left[\delta(k) \left\| \epsilon - \epsilon_{\theta} \left(\sqrt{\tilde{\alpha}_k} \mathbf{x}_t^0 + \sqrt{1 - \tilde{\alpha}_k} \epsilon, \mathbf{h}_{t-1}, k \right) \right\|^2 \right]. \quad (22)$$

The context window is used to generate the hidden state \mathbf{h}_{t_0-1} for the starting point of the training process. It is not hard to see that Equation (22) is very similar to Equation (6) except for the inclusion of hidden states to represent the historical information.

In the training process, the parameter θ is estimated by minimizing the negative log-likelihood objective function with stochastic sampling. Then, future time series are generated in a step-by-step manner. Suppose that the last time point of the complete time series is \tilde{T} . The first step is to derive the hidden state $\mathbf{h}_{\tilde{T}}$ based on the last available context window. Next, the observation for the next time point $\tilde{T} + 1$ is predicted in a similar way as DDPM:

$$\mathbf{x}_{\tilde{T}+1}^k \leftarrow \frac{\left(\mathbf{x}_{\tilde{T}+1}^{k+1} - \zeta(k+1) \epsilon_{\theta}(\mathbf{x}_{\tilde{T}+1}^{k+1}, \mathbf{h}_{\tilde{T}}, k+1) \right)}{\sqrt{\alpha_{k+1}}} + \sigma_{k+1} \mathbf{z},$$

The predicted $\mathbf{x}_{\tilde{T}+1}^k$ should be fed back to the RNN module to obtain $\mathbf{h}_{\tilde{T}+1}$ before the prediction for the next time point. The sampling process will be repeated until the desired length of the future time series is reached.

3.3 ScoreGrad

ScoreGrad shares the same target distribution as TimeGrad, but it is alternatively built upon SDEs, extending the diffusion process from discrete to continuous and replacing the number of diffusion steps with an interval of integration (Yan et al., 2021). ScoreGrad is composed of a feature extraction module and a conditional SDE-based score-matching module. The feature extraction module is almost identical to the computation of \mathbf{h}_t in TimeGrad. However, Yan et al. (2021) have discussed the potential of adopting other network structures to encode historical information, such as temporal convolutional networks (Oord et al., 2016) and attention-based networks (Vaswani et al., 2017). Here we still focus on RNN as the default choice. In the conditional SDE-based score matching module, the diffusion process is conducted through the same SDE as in Equation (9) but its associated time-reverse SDE is refined as following:

$$d\mathbf{x}_t = [f(\mathbf{x}_t, k) - g(k)^2 \nabla_{\mathbf{x}_t} \log q_k(\mathbf{x}_t | \mathbf{h}_t)] dk + g(k) d\mathbf{w}, \quad (23)$$

where $k \in [0, K]$ represents the SDE integral time. As a common practice, the conditional score function $\nabla_{\mathbf{x}_t} \log q_k(\mathbf{x}_t | \mathbf{h}_t)$ is approximated with a parameterized neural network $\mathbf{s}_\theta(\mathbf{x}_t^k, \mathbf{h}_t, k)$. Inspired by WaveNet (Oord et al., 2016) and DiffWave (Kong et al., 2021), the neural network is designed to have 8 connected residual blocks, while each block contains a bidirectional dilated convolution module, a gated activation unit, a skip-connection process, and an 1D convolutional neural network for output.

The objective function of ScoreGrad is a conditional modification of Equation (12), computed as

$$\sum_{t=t_0}^T L_t(\theta) \quad (24)$$

with $L_t(\theta)$ being

$$\mathbb{E}_{k, \mathbf{x}_t^0, \mathbf{x}_t^k} \left[\delta(k) \left\| \mathbf{s}_\theta(\mathbf{x}_t^k, \mathbf{h}_t, k) - \nabla_{\mathbf{x}_t} \log q_{0k}(\mathbf{x}_t | \mathbf{x}_t^0) \right\|^2 \right]. \quad (25)$$

Up to this point, we only use the general expression of SDE for simple illustration. In the training process, one shall decide the specific type of SDE to use. Potential options include VE SDE, VP SDE, and sub-VP SDE (Song et al., 2021). The optimization varies depending on the chosen SDE because different SDEs lead to different forward transition kernel $q(\mathbf{x}_t^k | \mathbf{x}_t^{k-1})$ and also different perturbation distribution $q_{0k}(\mathbf{x}_t | \mathbf{x}_t^0)$. Finally, for forecasting, ScoreGrad utilizes the predictor-corrector sampler as in (Song et al., 2021) to sample from the time-reverse SDE.

3.4 D³VAE

In practice, we may encounter the challenge of insufficient observations. If the historical multivariate time series were recorded based on a short period, they are prone to significant level of noises due to measurement errors, sampling variability, and randomness from other sources. To address the problem of limited and noisy time series, D³VAE,

proposed by Li et al. (2022), employs a coupled diffusion process for data augmentation, and then uses a bidirectional auto-encoder (BVAE) together with denoising score matching to clear the noise. In addition, D³VAE also considers disentangling latent variables by minimizing the overall correlation for better interpretability and stability of predictions. Moreover, the mean square error (MSE) between the prediction and actual observations in the prediction window is included in the objective function, further emphasizing the role of supervision.

Assuming that the prediction window can be generated from a set of latent variables \mathbf{Z} that follows a Gaussian distribution $q(\mathbf{Z}|\mathbf{x}_{1:t_0-1}^0)$. The conditional distribution of \mathbf{Z} is approximated with $p_\phi(\mathbf{Z}|\mathbf{x}_{1:t_0-1}^0)$ where ϕ denotes learnable parameters. Then, the forecasting time series $\hat{\mathbf{x}}_{t_0:T}$ can be generated from the estimated target distribution, given by $p_\theta(\mathbf{x}_{t_0:T}^0|\mathbf{Z})$. It is not difficult to see that the prediction window is still predicted based on the context window however with latent variables \mathbf{Z} as an intermediate.

In the coupled diffusion process, we inject noises separately into the context window and the prediction window. Different from TimeGrad which injects noises to the observation at each time point individually, the coupled diffusion process is applied to the whole period. For the context window, the same kernel as Equation (2) is applied such that

$$\mathbf{x}_{1:t_0-1}^k = \sqrt{\tilde{\alpha}_k} \mathbf{x}_{1:t_0-1}^0 + \sqrt{1 - \tilde{\alpha}_k} \boldsymbol{\epsilon}, \quad (26)$$

where $\boldsymbol{\epsilon}$ denotes the standard Gaussian noises but with a matrix rather than a vector form.

The diffusion process is further applied to the prediction window with adjusted noise levels $\alpha'_k > \alpha_k$. Let $\tilde{\alpha}'_k := \prod_{i=1}^k \alpha'_i$, then

$$\mathbf{x}_{t_0:T}^k = \sqrt{\tilde{\alpha}'_k} \mathbf{x}_{t_0:T}^0 + \sqrt{1 - \tilde{\alpha}'_k} \boldsymbol{\epsilon}. \quad (27)$$

This diffusion process simultaneously augments the context window and the prediction window, thus improving the generalization ability for short time series forecasting. Besides, it is proven by Li et al. (2022) that the uncertainty caused by the generative model and the inherent noises in the observed data can both be mitigated by the coupled diffusion process.

The backward process is accomplished with two steps. The first step is to predict $\mathbf{x}_{t_0:T}^k$ with a BVAE as the one used in (Vahdat and Kautz, 2020), which is composed of an encoder and a decoder with multiple residual blocks and takes the disturbed context window $\mathbf{x}_{1:t_0-1}^k$ as input. The latent variables in \mathbf{Z} are gradually generated and fed into the model in a summation manner. The output of this process is the predicted disturbed prediction window $\hat{\mathbf{x}}_{t_0:T}^k$. The second step involves further cleaning of the predicted data with a denoising score matching module. More specifically, the final prediction is obtained via a single-step gradient jump (Saremi and Hyvarinen, 2019):

$$\hat{\mathbf{x}}_{t_0:T}^0 \leftarrow \hat{\mathbf{x}}_{t_0:T}^k - \sigma_0^2 \nabla_{\hat{\mathbf{x}}_{t_0:T}^k} E(\hat{\mathbf{x}}_{t_0:T}^k; e),$$

where σ_0 is prescribed and $E(\hat{\mathbf{x}}_{t_0:T}^k; e)$ is the energy function.

Disentanglement of latent variables \mathbf{Z} can efficiently enhance the model interpretability and reliability for prediction (Li et al., 2021). It is measured by the total correlation

of the random latent variables Z . Generally, a lower total correlation implies better disentanglement which is a signal of useful information. The computation of total correlation happens synchronously with the BVAE module.

The objective function of D³VAE consists of four components. It can be written as

$$w_1 D_{KL} \left(q(\mathbf{x}_{t_0:T}^k) \| p_{\theta}(\hat{\mathbf{x}}_{t_0:T}^k) \right) + w_2 \mathcal{L}_{DSM} + w_3 \mathcal{L}_{TC} + \mathcal{L}_{MSE}, \quad (28)$$

where w_1, w_2, w_3 are trade-off parameters that assign significance levels to the components. The first component $D_{KL} \left(q(\mathbf{x}_{t_0:T}^k) \| p_{\theta}(\hat{\mathbf{x}}_{t_0:T}^k) \right)$ matches the estimated target distribution with the true distribution of the prediction window. The last three components, $\mathcal{L}_{DSM}, \mathcal{L}_{TC}$, and \mathcal{L}_{MSE} , are corresponding to the DSM module, disentanglement of latent variables, and MSE between the prediction and the truth, respectively. The parameters θ and ϕ are learned together in the training process. Eventually, forecasting samples are generated from the learned distribution $p_{\phi}(\mathbf{Z} | \mathbf{x}_{1:t_0-1}^0)$ and $p_{\theta}(\mathbf{x}_{t_0:T}^0 | \mathbf{Z})$.

3.5 DSPD

Multivariate time series data can be considered as a record of value changes for multiple features of an entity of interest. Data are collected from the same entity, and the measuring tools normally stay unchanged during the whole observed time period. So, assuming that the change of variables over time is smooth, the time series data can be modelled as values from an underlying continuous function (Biloš et al., 2022). In this case, the context window is expressed as $\mathbf{X}_c^0 = \{\mathbf{x}(1), \mathbf{x}(2), \dots, \mathbf{x}(t_0 - 1)\}$ and the prediction window becomes $\mathbf{X}_p^0 = \{\mathbf{x}(t_0), \mathbf{x}(t_0 + 1), \dots, \mathbf{x}(T)\}$, where $\mathbf{x}(\cdot)$ is a continuous function of the time point t .

Different from traditional diffusion models, the diffusion and reverse processes are no longer applied to vector observations at each time point. Alternatively, the target of interest is the continuous function $\mathbf{x}(\cdot)$, which means noises will be injected and removed from a function rather than a vector. Therefore, a continuous noise function $\epsilon(\cdot)$ should take place of the noise vector $\epsilon \sim \mathcal{N}(\mathbf{0}, \mathbf{I})$. This function should be both continuous and tractable such that it accounts for the correlation between measurements and enables training and sampling. These requirements are effectively satisfied by designing a Gaussian stochastic process $\epsilon(\cdot) \sim \mathcal{GP}(\mathbf{0}, \mathbf{\Sigma})$ (Biloš et al., 2022).

Discrete stochastic process diffusion (DSPD) is built upon the DDPM formulation but with the stochastic process $\epsilon(\cdot) \sim \mathcal{GP}(\mathbf{0}, \mathbf{\Sigma})$. It is a delight that DSPD is only slightly different from DDPM in terms of implementation. More specifically, DSPD simply replaces the commonly applied noise $\epsilon \sim \mathcal{N}(\mathbf{0}, \mathbf{I})$ with the noise function $\epsilon(\cdot)$ whose discretized form is $\epsilon \sim \mathcal{N}(\mathbf{0}, \mathbf{\Sigma})$. Let \mathbf{X}^0 be an observed multivariate time series in a certain period of time $\mathbf{T}' = \{t'_1, t'_2, \dots, t'_T\}$, which means $\mathbf{X}^0 = \{\mathbf{x}(t'_1), \mathbf{x}(t'_2), \dots, \mathbf{x}(t'_T)\}$. In the forward process, noises are injected through the transition kernel

$$q(\mathbf{X}^k | \mathbf{X}^0) = \mathcal{N}(\sqrt{\tilde{\alpha}_k} \mathbf{X}^0, (1 - \tilde{\alpha}_k) \mathbf{\Sigma}) \quad (29)$$

Then, the following backward transition kernel is applied to recover the original data:

$$p_{\theta}(\mathbf{X}^{k-1} | \mathbf{X}^k) = \mathcal{N}(\mu_{\theta}(\mathbf{X}^k, k), (1 - \alpha_k) \mathbf{\Sigma}). \quad (30)$$

Consequently, the objective function should be changed as

$$\mathbb{E}_{k, \mathbf{X}^0, \boldsymbol{\epsilon}} \left[\delta(k) \left\| \boldsymbol{\epsilon} - \boldsymbol{\epsilon}_\theta \left(\sqrt{\tilde{\alpha}_k} \mathbf{X}^0 + \sqrt{1 - \tilde{\alpha}_k} \boldsymbol{\epsilon}, k \right) \right\|^2 \right], \quad (31)$$

where $\boldsymbol{\epsilon} \sim \mathcal{N}(\mathbf{0}, \boldsymbol{\Sigma})$ with the covariance matrix from the Gaussian process $\mathcal{GP}(\mathbf{0}, \boldsymbol{\Sigma})$.

Forecasting via DSPD is very similar to TimeGrad. As before, the aim is still to learn the conditional probability $q(\mathbf{X}_p^0 | \mathbf{X}_c^0)$. But there are two major improvements. Firstly, the prediction is available for any future time point in the continuous time interval. Secondly, instead of step-by-step forecasting, DSPD can generate samples for multiple time points in one run. By adding the historical condition into the fundamental objective function in Equation (31), the objective function for DSPD forecasting is then given by

$$\mathbb{E}_{k, \mathbf{X}_p^0, \boldsymbol{\epsilon}} \left[\delta(k) \left\| \boldsymbol{\epsilon} - \boldsymbol{\epsilon}_\theta \left(\sqrt{\tilde{\alpha}_k} \mathbf{X}_p^0 + \sqrt{1 - \tilde{\alpha}_k} \boldsymbol{\epsilon}, \mathbf{X}_c^0, k \right) \right\|^2 \right]. \quad (32)$$

Finally, supposing that the last context window is $\tilde{\mathbf{X}}_c$, the sampling process to forecast for the prediction target $\tilde{\mathbf{X}}_p$ in time interval $\tilde{\mathbf{T}}$ is given by

$$\tilde{\mathbf{X}}_p^k \leftarrow \frac{\left(\tilde{\mathbf{X}}_p^{k+1} - \zeta(k+1) \mathbf{L} \boldsymbol{\epsilon}_\theta(\tilde{\mathbf{X}}_p^{k+1}, \tilde{\mathbf{X}}_c, k+1) \right)}{\sqrt{\alpha_{k+1}}} + (1 - \alpha_{k+1}) \mathbf{z},$$

where \mathbf{L} is from the factorization of the covariance matrix $\boldsymbol{\Sigma} = \mathbf{L} \mathbf{L}^\top$, the last diffusion output is generated as $\tilde{\mathbf{X}}_p^K \sim \mathcal{N}(\mathbf{0}, \boldsymbol{\Sigma})$, and $\mathbf{z} \sim \mathcal{N}(\mathbf{0}, \boldsymbol{\Sigma})$.

Similar to the extension from TimeGrad to ScoreGrad, the continuous noise function can also be adapted into the SDE framework, thus leading to the continuous stochastic process diffusion (CSPD) model (Biloš et al., 2022). The diffusion process of CSPD introduces the factorized covariance matrix $\boldsymbol{\Sigma} = \mathbf{L} \mathbf{L}^\top$ to VP SDE (see section 2.3) as

$$d\mathbf{X} = -\frac{1}{2} \alpha(k) \mathbf{X} dk + \sqrt{\alpha(k)} \mathbf{L} d\mathbf{w}, \quad (33)$$

where \mathbf{w} is a matrix that represents a standard Wiener process. The perturbation distribution in the objective function is then modified as

$$q_{0k}(\mathbf{X}^k | \mathbf{X}^0) = \mathcal{N} \left(\mathbf{X}^k; \mathbf{X}^0 e^{-\frac{1}{2} \int_0^k \alpha(s) ds}, [1 - e^{-\int_0^k \alpha(s) ds}] \boldsymbol{\Sigma} \right). \quad (34)$$

3.6 DiffSTG

Spatio-temporal graphs (STGs) are a special type of multivariate time series that encodes spatial and temporal relationships and interactions among different entities in a graph structure (Wen et al., 2023). They are commonly observed in real-life applications such as traffic flow prediction (Li et al., 2018), weather forecasting (Simeunović et al., 2021), and finance prediction (Zhou et al., 2011). Suppose we have N entities of interest, such as traffic sensors or companies in the stock market. We can model these entities and their underlying relationships as a graph $\mathcal{G} = \{\mathcal{V}, \mathcal{E}, \mathbf{W}\}$, where \mathcal{V} is a set of N nodes as representations for entities, \mathcal{E} is a set of links that indicates the relationship between

nodes, and \mathbf{W} is a weighted adjacency matrix that describes the graph topological structure. Multivariate time series observed at all entities are models as graph signals $\mathbf{X}_c^0 = \{\mathbf{x}_1^0, \mathbf{x}_2^0, \dots, \mathbf{x}_{t_0-1}^0 | \mathbf{x}_t^0 \in \mathbb{R}^{D \times N}\}$, which means we have D -dimensional observations from N entities at each time point t . Identical to the previous problem formulation, the aim of STG forecasting is also to predict $\mathbf{X}_p^0 = \{\mathbf{x}_{t_0}^0, \mathbf{x}_{t_0+1}^0, \dots, \mathbf{x}_T^0 | \mathbf{x}_t^0 \in \mathbb{R}^{D \times N}\}$ based on the historical information \mathbf{X}_c . Nevertheless, except for the time dependency on historical observations, we also need to consider the spatial interactions between different entities represented by the graph topology.

DiffSTG applies diffusion models on STG forecasting with a graph-based noise-matching network called UGnet (Wen et al., 2023). The idea of DiffSTG can be regarded as the extension of DDPM-based forecasting to STGs with an additional condition on the graph structure, which means the target distribution in Equation (17) is approximated alternatively by

$$p_{\theta}(\mathbf{x}_{t_0:T}^0 | \mathbf{x}_{1:t_0-1}^0, \mathbf{W}). \quad (35)$$

Accordingly, the objective function is changed as

$$\mathbb{E}_{k, \mathbf{x}_{t_0:T}^0, \epsilon} \left[\delta(k) \left\| \epsilon - \epsilon_{\theta} \left(\sqrt{\tilde{\alpha}_k} \mathbf{x}_{t_0:T}^0 + \sqrt{1 - \tilde{\alpha}_k} \epsilon, \mathbf{x}_{1:t_0-1}^0, k, \mathbf{W} \right) \right\|^2 \right]. \quad (36)$$

The objective function in Equation (36) actually treats the context window and the prediction window as samples from two separate sample spaces, namely, $\mathbf{X}_c^0 \in \mathcal{X}_c$ and $\mathbf{X}_p^0 \in \mathcal{X}_p$ with \mathcal{X}_c and \mathcal{X}_p being two individual sample spaces. However, considering the fact that the context and prediction intervals are consecutive, it may be more reasonable to treat the two windows as a complete sample from the same sample space. To this end, Wen et al. (2023) reformulate the forecasting problem and revise the approximation in Equation (35) as

$$p_{\theta}(\mathbf{x}_{1:T}^0 | \mathbf{x}_{1:t_0-1}^0, \mathbf{W}), \quad (37)$$

in which the history condition is derived by masking the future time series from the whole time period. The associated objective function is

$$\mathbb{E}_{k, \mathbf{x}_{1:T}^0, \epsilon} \left[\delta(k) \left\| \epsilon - \epsilon_{\theta} \left(\sqrt{\tilde{\alpha}_k} \mathbf{x}_{1:T}^0 + \sqrt{1 - \tilde{\alpha}_k} \epsilon, \mathbf{x}_{1:t_0-1}^0, k, \mathbf{W} \right) \right\|^2 \right]. \quad (38)$$

The training process is quite straightforward following the common practice. But it should be noted that the sample generated in the forecasting process includes both historical and future values. So, we need to take out the forecasting target in the sample as the prediction.

Now, there is only one remaining problem. How to encode the graph structural information in the noise-matching network ϵ_{θ} ? Wen et al. (2023) proposed UGnet, an Unet-based network architecture (Ronneberger et al., 2015) combined with a graph neural network (GNN) to process time dependency and spatial relationships simultaneously. UGnet takes $\mathbf{x}_{1:T}^k, \mathbf{x}_{1:t_0-1}^0, k$ and \mathbf{W} as inputs and then outputs the prediction of the associated error ϵ .

3.7 GCRDD

Graph convolutional recurrent denoising diffusion model (GCRDD) is another diffusion-based model for STG forecasting (Li et al., 2023). It differs from DiffSTG that it uses hidden states from a recurrent component to store historical information as TimeGrad and employs a different network structure for the noise-matching term ϵ_θ . Please note that the notations related to STGs here follow subsection 3.6.

GCRDD approximates the target distribution with a probabilistic density function conditional on the hidden states and graph structure as following:

$$\prod_{t=\ell_0}^T p_\theta(\mathbf{x}_t^0 | \mathbf{h}_{t-1}, \mathbf{W}), \quad (39)$$

where the hidden state is computed with a graph-modified GRU, written as

$$\mathbf{h}_t = \text{GraphGRU}_\theta(\mathbf{x}_t^0, \mathbf{h}_{t-1}, \mathbf{W}). \quad (40)$$

The graph-modified GRU replaces the weight matrix multiplication in traditional GRU (Chung et al., 2014) with graph convolution such that both temporal and spatial information is stored in the hidden state. The objective function of GCRDD adopts a similar form of TimeGrad but with additional graph structural information in the noise-matching network:

$$\mathbb{E}_{k, \mathbf{x}_t^0, \epsilon} \left[\delta(k) \left\| \epsilon - \epsilon_\theta \left(\sqrt{\tilde{\alpha}_k} \mathbf{x}_t^0 + \sqrt{1 - \tilde{\alpha}_k} \epsilon, \mathbf{h}_{t-1}, \mathbf{W}, k \right) \right\|^2 \right]. \quad (41)$$

For the noise-matching term, GCRDD adopts a variant of DiffWave (Kong et al., 2021) that incorporates a graph convolution component to process spatial information in \mathbf{W} . STG forecasting via GCRDD is the same as TimeGrad except that the sample generated at each time point is a matrix rather than a vector.

4 Time Series Imputation

In real-world problem settings, we usually encounter the challenge of missing values. When collecting time series data, the collection conditions may change over time, which makes it difficult to ensure the completeness of observation. In addition, accidents such as sensor failures and human errors may also result in the missing of historical records. Missing values in time series data normally have a negative impact on the accuracy of analysis and forecasting since the lack of partial observations makes the inference and conclusions vulnerable in future generalization.

Time series imputation aims to fill in the missing values in incomplete time series data. Many previous studies have focused on designing deep learning-based algorithms for time series imputation (Osman et al., 2018). Most existing approaches involve the RNN architecture to encode time-dependency in the imputation task (Che et al., 2018; Cao et al., 2018; Luo et al., 2018; Yoon et al., 2018). Except for these deterministic methods, probabilistic imputation models such as GP-VAE (Fortuin et al., 2020) and V-RIN (Mulyadi et al., 2021) have also shown their practical value in recent years. As

a rising star in probabilistic models, diffusion models have also been applied to time series imputation tasks (Tashiro et al., 2021; Alcaraz and Strodthoff, 2023). Compared with other probabilistic approaches, diffusion-based imputation enjoys high flexibility in the assumption of the true data distribution. In this section, we will cover four diffusion-based methods, including three for multivariate time series imputation and one for STG imputation.

4.1 Problem Formulation

We still consider the multivariate time series $\mathbf{X}^0 = \{\mathbf{x}_1^0, \mathbf{x}_2^0, \dots, \mathbf{x}_T^0 | \mathbf{x}_i^0 \in \mathbb{R}^D\}$. It is not difficult to see that $\mathbf{X}^0 \in \mathbb{R}^{D \times T}$, where D is the number of features, and T is the number of time points in the period $[1, T]$. Different from time series forecasting in which we assume that all elements in \mathbf{X}^0 are known, here we have an incomplete matrix of observations. In the imputation task, we try to predict the values of missing data by exploring the information from some observed data. We denote the observed data as \mathbf{X}_{ob}^0 and the missing data as \mathbf{X}_{ms}^0 . Then, the imputation task is to find the conditional probability distribution $q(\mathbf{X}_{ms}^0 | \mathbf{X}_{ob}^0)$.

For practical purposes, zero padding is applied to the incomplete matrix \mathbf{X}^0 such that all missing entries are assigned to 0. In addition, a zero-one matrix $\mathbf{M} \in \mathbb{R}^{D \times T}$ is constructed as a mask to denote the position of missing values. More specifically, the elements in \mathbf{M} is 0 when the corresponding value in \mathbf{X}^0 is missing, and 1 otherwise.

Both \mathbf{X}_{ob}^0 and \mathbf{X}_{ms}^0 have the same dimension as \mathbf{X}^0 . In the training process, a fraction of the actually observed data in \mathbf{X}^0 is randomly selected to be the true values of missing data, and the rest of the observed data will be the condition for prediction. A training mask $\mathbf{M}' \in \mathbb{R}^{D \times T}$ is introduced to obtain \mathbf{X}_{ob}^0 and \mathbf{X}_{ms}^0 . It is constructed by assigning 1 to entries that are corresponding to the remaining observed data in \mathbf{X}^0 . Then, \mathbf{X}_{ob}^0 is computed as $\mathbf{X}_{ob}^0 = \mathbf{M}' \odot \mathbf{X}^0$, and \mathbf{X}_{ms}^0 is computed as $\mathbf{X}_{ms}^0 = (\mathbf{M} - \mathbf{M}') \odot \mathbf{X}^0$, where \odot denotes the element-wise matrix multiplication. In the forecasting process, on the other hand, all actually observed data are used as the condition, which means $\mathbf{X}_{ob}^0 = \mathbf{M} \odot \mathbf{X}^0$.

It is worth mentioning that the problem formulation here is only a typical case. We will introduce later in subsection 4.4 about another formulation that takes the whole time series matrix \mathbf{X}^0 as the target for generation.

4.2 CSDI

Conditional Score-based Diffusion model for Imputation (CSDI) is the pioneering work on diffusion-based time series imputation (Tashiro et al., 2021). Identical to TimeGrad, the basic diffusion formulation of CSDI is also DDPM. However, as we have discussed in section 3.2, the historical information is encoded by an RNN module in TimeGrad, which hampers the direct extension of TimeGrad to imputation tasks because the computation of hidden states may be interrupted by missing values in the context window.

CSDI applies the diffusion and reverse processes to the matrix of missing data, \mathbf{X}_{ms}^0 . Correspondingly, the reverse transition kernel is refined as a probabilistic distribution

conditional on \mathbf{X}_{ob}^0 :

$$\begin{aligned} p_{\theta}(\mathbf{X}_{ms}^{k-1} | \mathbf{X}_{ms}^k, \mathbf{X}_{ob}^0) \\ = \mathcal{N}(\mathbf{X}_{ms}^{k-1}; \boldsymbol{\mu}_{\theta}(\mathbf{X}_{ms}^k, k | \mathbf{X}_{ob}^0), \sigma_{\theta}(\mathbf{X}_{ms}^k, k | \mathbf{X}_{ob}^0) \mathbf{I}), \end{aligned} \quad (42)$$

where

$$\boldsymbol{\mu}_{\theta}(\mathbf{X}_{ms}^k, k | \mathbf{X}_{ob}^0) = \frac{1}{\sqrt{\alpha_k}} \left(\mathbf{X}_{ms}^k - \zeta(k) \boldsymbol{\epsilon}_{\theta}(\mathbf{X}_{ms}^k, k | \mathbf{X}_{ob}^0) \right) \quad (43)$$

with $\mathbf{X}_{ms}^k = \sqrt{\alpha_k} \mathbf{X}_{ms}^0 + \sqrt{1 - \alpha_k} \boldsymbol{\epsilon}$. One may notice that the variance term here is different from the version in DDPM (Ho et al., 2020). Previously, the variance term is defined with some pre-specified constant σ_k with $k = 1, 2, \dots, K$, implying that the variance is treated as a hyperparameter. CSDI, however, defines a learnable version σ_{θ} with parameter θ . Both ways are acceptable and have their respective practical value.

The objective function of CSDI is given by

$$\mathbb{E}_{k, \mathbf{X}_{ms}^0, \boldsymbol{\epsilon}} \left[\delta(k) \left\| \boldsymbol{\epsilon} - \boldsymbol{\epsilon}_{\theta} \left(\sqrt{\alpha_k} \mathbf{X}_{ms}^0 + \sqrt{1 - \alpha_k} \boldsymbol{\epsilon}, \mathbf{X}_{ob}^0, k \right) \right\|^2 \right]. \quad (44)$$

The noise-matching network $\boldsymbol{\epsilon}_{\theta}$ adopts the DiffWave (Kong et al., 2021) architecture by default. After training, the imputation is accomplished by generating the target matrix of missing values in the same way as DDPM. \mathbf{X}_{ob}^0 in the sampling process is identical to the zero padding version of the original time series matrix \mathbf{X}^0 , where all missing values are assigned to 0. The starting point of the sampling process is a random Gaussian imputation target $\mathbf{X}_{ms}^K \sim \mathcal{N}(\mathbf{0}, \mathbf{I})$. Then, for $k = K - 1, \dots, 1$, the algorithm computes:

$$\mathbf{X}_{ms}^k \leftarrow \frac{(\mathbf{X}_{ms}^{k+1} - \zeta(k+1) \boldsymbol{\epsilon}_{\theta}(\mathbf{X}_{ms}^{k+1}, \mathbf{X}_{ob}^0, k+1))}{\sqrt{\alpha_{k+1}}} + \sigma_{\theta} \mathbf{Z},$$

where $\mathbf{Z} \sim \mathcal{N}(\mathbf{0}, \mathbf{I})$ for $k = K - 1, \dots, 1$, and $\mathbf{Z} = \mathbf{0}$ for $k = 0$.

4.3 DSPD

Here we discuss the simple extension of DSPD and CSPD in section 3.5 to time series imputation tasks. Since the rationale behind the extension of these two models is almost identical, we only focus on DSPD for illustration. The assumption of DSPD states that the observed time series is formed by values of a continuous function $\boldsymbol{x}(\cdot)$ of time t . Therefore, the missing values can be obtained by computing the values of this continuous function at the corresponding time points. Recall that DSPD utilizes the covariance matrix $\boldsymbol{\Sigma}$ instead of the DDPM variance $\sigma_k \mathbf{I}$ or σ_{θ} in the backward process. Therefore, one may apply DSPD to imputation tasks in a similar way as CSDI by replacing the variance term in Equation (42) with the covariance matrix $\boldsymbol{\Sigma}$. According to Biloš et al. (2022), the continuous noise process is a more natural choice than the discrete noise vector because it takes account of the irregularity in the measurement when collecting the time series data.

4.4 SSSD

Structured state space diffusion (SSSD) differs from the aforementioned two methods by having the whole time series matrix \mathbf{X}^0 as the generative target in its diffusion module (Alcaraz and Strodthoff, 2023). The name, “structured state space diffusion”, comes from the design of the noise-matching network ϵ_θ , which adopts the state space model (Gu et al., 2022) as the internal architecture. As a matter of fact, ϵ_θ can also take other architectures such as the DiffWave-based network in CSDI (Tashiro et al., 2021) and SaShiMi, a generative model for sequential data (Goel et al., 2022). However, the authors of SSSD have shown empirically that the structured state space model generally generates the best imputation outcome compared with other architectures (Alcaraz and Strodthoff, 2023). To emphasize the difference between this method with other diffusion-based approaches, here we will primarily focus on the unique problem formulation used by SSSD.

As we have mentioned, the generative target of SSSD is the whole time series matrix, $\mathbf{X}^0 \in \mathbb{R}^{D \times T}$, rather than a matrix that particularly represents the missing values. For the purpose of training, \mathbf{X}^0 is also processed with zero padding. The conditional information, in this case, is from a concatenated matrix $\mathbf{X}_c^0 = \text{Concat}(\mathbf{X}^0 \odot \mathbf{M}_c, \mathbf{M}_c)$, where \mathbf{M}_c is a zero-one matrix indicating the position of observed values as the condition. The element in \mathbf{M}_c can only be 1 if its corresponding value in \mathbf{X}^0 is known.

There are two options for the objective function used in the training process. Similar to other approaches, the objective function can be a simple conditional variant of the DDPM objective function:

$$\mathbb{E}_{k, \mathbf{X}^0, \epsilon} \left[\delta(k) \left\| \epsilon - \epsilon_\theta \left(\sqrt{\tilde{\alpha}_k} \mathbf{X}^0 + \sqrt{1 - \tilde{\alpha}_k} \epsilon, \mathbf{X}_c^0, k \right) \right\|^2 \right], \quad (45)$$

where ϵ_θ is by default built upon the structured state space model. The other choice of the objective function is computed with only known data, which is mathematically expressed as

$$\mathbb{E}_{k, \mathbf{X}^0, \epsilon} \left[\delta(k) \left\| \epsilon \odot \mathbf{M}_c - \epsilon_\theta \left(\sqrt{\tilde{\alpha}_k} \mathbf{X}^0 + \sqrt{1 - \tilde{\alpha}_k} \epsilon, \mathbf{X}_c^0, k \right) \odot \mathbf{M}_c \right\|^2 \right]. \quad (46)$$

According to Alcaraz and Strodthoff (2023), the second objective function is typically a better choice in practice. For forecasting, SSSD employs the usual sampling algorithm and applies to the unknown entries in \mathbf{X}^0 , namely, $(1 - \mathbf{M}_c) \odot \mathbf{X}^0$.

An interesting point proposed along with SSSD is that imputation models can also be applied to forecasting tasks. This is because future time series can be viewed as a long block of missing values on the right of \mathbf{X}^0 . Nevertheless, experiments of SSSD have shown that the diffusion-based approaches underperform other methods such as the Autoformer (Wu et al., 2021) in forecasting tasks.

4.5 PriSTI

PriSTI is a diffusion-based model for STG imputation (Liu et al., 2023). However, different from DiffSTG, the existing framework of PriSTI is designed for STGs with only

one feature, which means the graph signal has the form $\mathbf{X}^0 = \{\mathbf{x}_1^0, \mathbf{x}_2^0, \dots, \mathbf{x}_T^0\} \in \mathbb{R}^{N \times T}$. Each vector $\mathbf{x}_t^0 \in \mathbb{R}^N$ represents the observed values of N nodes at time point t . This kind of data is often observed in traffic prediction (Li et al., 2018) and weather forecasting (Yi et al., 2016). *METR-LA*, for example, is an STG dataset that contains traffic speed collected by 207 sensors on a Los Angeles highway in a 4-month time period (Li et al., 2018). There is only one node attribute, that is, traffic speed. However, unlike the multivariate time series matrix, where features (sensors in this case) are usually assumed to be uncorrelated, the geographic relationship between different sensors is stored in the weighted adjacency matrix \mathbf{W} , allowing a more pertinent representation of real-world traffic data.

The number of nodes N can be considered as the number of features D in CSDI. The only difference is that PriSTI incorporates the underlying relationship between each pair of nodes in the conditional information for imputation. So, the problem formulation adopted by PriSTI is the same as our discussion in subsection 4.1, thus the goal is still to find $q(\mathbf{X}_{ms}^0 | \mathbf{X}_{ob}^0)$.

To encode graph structural information, the mean in Equation (43) is modified as

$$\boldsymbol{\mu}_\theta(\mathbf{X}_{ms}^k, k | \mathbf{X}_{ob}^0, \mathbf{W}) = \frac{1}{\sqrt{\alpha_k}} \left(\mathbf{X}_{ms}^k - \zeta(k) \boldsymbol{\epsilon}_\theta(\mathbf{X}_{ms}^k, \mathbf{X}_{ob}^0, k, \mathbf{W}) \right), \quad (47)$$

where \mathbf{W} is the weighted adjacency matrix. Consequently, the objective function is changed as

$$\mathbb{E}_{k, \mathbf{X}_{ms}^0, \boldsymbol{\epsilon}} \left[\left\| \delta(k) \left\| \boldsymbol{\epsilon} - \boldsymbol{\epsilon}_\theta \left(\sqrt{\tilde{\alpha}_k} \mathbf{X}_{ms}^0 + \sqrt{1 - \tilde{\alpha}_k} \boldsymbol{\epsilon}, \mathbf{X}_{ob}^0, k, \mathbf{W} \right) \right\|^2 \right]. \quad (48)$$

The conditional information, \mathbf{X}_{ob}^0 , is processed with linear interpolation before it is fed into the algorithm to incorporate extra noises, which will enhance the denoising capability of the model and eventually lead to better consistency in prediction (Choi et al., 2022). The noise-matching network $\boldsymbol{\epsilon}_\theta$ is composed of two modules, including a conditional feature extraction module and a noise estimation module. The conditional feature extraction module takes the interpolated information \mathbf{X}_{ob}^0 and adjacency matrix \mathbf{W} as inputs and generates a global context with both spatial and temporal information as the condition for diffusion. Then, the noise estimation module utilizes this global context to estimate the injected noises with a specialized attention mechanism to capture temporal dependencies and geographic information. Ultimately, the STG imputation is fulfilled with the usual sampling process of DDPM, but with the specially designed noise-matching network here to incorporate the additional spatial relationship.

Since PriSTI only works for the imputation of STGs with a single feature, which is simply a special case of STGs, this model’s practical value is somehow limited. So, the extension of the idea here to more generalized STGs is a notable topic for future researchers.

5 Time Series Generation

The rapid development of the machine learning paradigm requires high-quality data for different learning tasks in finance, economics, physics, and other fields. The performance

of machine learning model and algorithm may highly subject to the underlying data quality. Time series generation refers to the process of creating synthetic data that resembles the real-world time series. Since the time series data is characterized by its temporal dependencies, the generation process usually requires the learning of underlying patterns and trends, from the past information.

Time series generation is a developing topic in the literature with the existence of several methods (Yoon et al., 2018; Desai et al., 2021). Time series data can be seen as a case of sequential data, whose generation usually involves the GAN architecture (Xu et al., 2020; Donahue et al., 2018; Esteban et al., 2017; Mogren, 2016). Accordingly, TimeGAN is proposed to generate time series data based on an integration of RNN and GAN for the purpose of processing time dependency and generation (Yoon et al., 2018). However, the GAN-based generative methods have been criticized as they are unstable (Chu et al., 2020) and subject to the model collapse issue (Xiao et al., 2021). Another way to generate time series data is stemmed from the variational autoencoder, leading to the so-called TimeVAE model (Desai et al., 2021). As a common shortcoming of VAE-based models, TimeVAE requires a user-defined distribution for its probabilistic process. Here we will present a different probabilistic time series generator originated from diffusion models, which is more flexible with the form of the target distribution. We will particularly focus on (Lim et al., 2023) because it is the first and only work on this novel design. This section aims to enlighten researchers about this new-born research direction, and we expect to see more derivative works in the future.

5.1 Problem Formulation

With the multivariate time series $\mathbf{X}^0 = \{\mathbf{x}_1^0, \mathbf{x}_2^0, \dots, \mathbf{x}_T^0 | \mathbf{x}_i^0 \in \mathbb{R}^D\}$, the time series generation problem aims to synthesize time series $\mathbf{x}_{1:T}^0$ by generating observation \mathbf{x}_t^0 at time point $t \in [2, T]$ with the consideration of its previous historical data $\mathbf{x}_{1:t-1}^0$. Correspondingly, the target distribution is the conditional density $q(\mathbf{x}_t^0 | \mathbf{x}_{1:t-1}^0)$ for $t \in [2, T]$, and the associated generative process involves the recursive sampling of \mathbf{x}_t for all time points in the observed period. Details about the training and generation processes will be discussed in the next subsection.

5.2 TSGM

To our best knowledge, (Lim et al., 2023) is the only work to study the time series generation problem based on the diffusion method. The conditional score-based time series generative model (TSGM) was proposed, to conditionally generate each time series observation based on the past generated observations. The TSGM architecture includes three components: an encoder, a decoder and a conditional score-matching network. The pre-trained encoder is used to embed the underlying time series into a latent space. The conditional score-matching network is used to sample the hidden states, which are then converted to the time series samples via the decoder.

Given the multivariate time series $\mathbf{X}^0 = \{\mathbf{x}_1^0, \mathbf{x}_2^0, \dots, \mathbf{x}_T^0 | \mathbf{x}_i^0 \in \mathbb{R}^D\}$, the encoder **En** and decoder **De** enable the mapping between the time series data and hidden states in

a latent space. The mapping process can be defined as:

$$\mathbf{h}_t = \mathbf{En}(\mathbf{h}_{t-1}^0, \mathbf{x}_t^0), \quad \hat{\mathbf{x}}_t^0 = \mathbf{De}(\mathbf{h}_t^0), \quad (49)$$

where $\hat{\mathbf{x}}_t^0$ refers to the reconstructed time series data at time t , after the mapping process. This is a recursive process as both the encoder \mathbf{En} and decoder \mathbf{De} are constructed with the RNN structure. The training objective function \mathcal{L}_{ED} for both encoder and decoder is defined as:

$$\mathcal{L}_{ED} = \mathbb{E}_{\mathbf{x}_{1:T}^0} [\|\hat{\mathbf{x}}_{1:T}^0 - \mathbf{x}_{1:T}^0\|_2^2]. \quad (50)$$

Given the auto-dependency characteristic of the time series data, learning the conditional log-likelihood function is essential. To address this, the conditional score-matching network is designed based on the SDE formulation of diffusion models. It is worth noting that TSGM focuses on the generation of hidden states rather than producing the time series directly with the sampling process. At time step t , instead of applying the diffusion process to \mathbf{x}_t^0 , the hidden states \mathbf{h}_t^0 is diffused to a Gaussian distribution by the following forward SDE:

$$d\mathbf{h}_t = f(k, \mathbf{h}_t)dk + g(k)d\boldsymbol{\omega} \quad (51)$$

where $k \in [0, K]$ refers to the integral time. With the diffused sample $\mathbf{h}_{1:t}^k$, the conditional score-matching network s_θ learns the gradient of the conditional log-likelihood function with the following objective function:

$$\mathcal{L}_{Score} = \mathbb{E}_{\mathbf{h}_{1:T}^0, \mathbf{k}} \sum_{t=1}^T [\mathcal{L}(t, k)], \quad (52)$$

with

$$\mathcal{L}(t, k) = \mathbb{E}_{\mathbf{h}_t^k} \left[\delta(k) \|s_\theta(\mathbf{h}_t^k, \mathbf{h}_{t-1}, k) - \nabla_{\mathbf{h}_t} \log q_{0k}(\mathbf{h}_t | \mathbf{h}_t^0)\|^2 \right]. \quad (53)$$

The network architecture of s_θ is designed based on U-net (Ronneberger et al., 2015), which was adopted by the classic SDE model (Song et al., 2021).

In the training process, the encoder and decoder are pre-trained using the objective \mathcal{L}_{ED} . They can also be trained simultaneously with the network s_θ , but Lim et al. (2023) showed that the pre-training generally led to better performance. Then, to learn the score-matching network, hidden states are firstly obtained through inputting the entire time series $\mathbf{x}_{1:T}^0$ into the encoder, and then fed into the training algorithm with the objective function \mathcal{L}_{Score} . The time series generation is achieved by sampling hidden states and then applying the decoder, where the sampling process is analogous to solving the solutions to the time-reverse SDE.

The TSGM method can achieve state-of-the-art sampling quality and diversity, compared to a range of well-developed time series generation methods. However, it is still subject to the fundamental limitation that all diffusion models may have: they are generally more computationally expensive than GANs.

6 Conclusion

Diffusion models, a rising star in advanced generative techniques, have shown their exceptional power in various real-world applications. In recent years, many successful attempts have been made to incorporate diffusion in time series applications to boost model performance. As a compensation for the deficiency of a methodical summary and discourse on the diffusion-based approaches for time series, we have furnished a self-contained survey of these approaches while discussing the interactions and differences among them. More specifically, we have presented six models for time series forecasting, four models for time series imputation, and one model for time series generation. Although these models have shown good performance with empirical evidence, we feel obligated to emphasize that they are usually associated with very high computational costs. In addition, since most models are constructed with a high level of theoretical background, there is still a lack of deeper discussion and exploration of the rationale behind these models. This survey is expected to serve as a starting point for new researchers in this area and also an inspiration for future directions.

References

- Alcaraz JL, Strodthoff N, 2023. Diffusion-based time series imputation and forecasting with structured state space models. *Transactions on Machine Learning Research*, 3.
- Anderson BD, 1982. Reverse-time diffusion equation models. *Stochastic Processes and their Applications*, 12(3):313-326.
- Austin J, Johnson DD, Ho J, et al., 2021. Structured denoising diffusion models in discrete state-spaces. *Advances in Neural Information Processing Systems*, 34:17981-17993.
- Biloš M, Rasul K, Schneider A, et al., 2022. Modeling temporal data as continuous functions with process diffusion. *arXiv preprint arXiv:221102590*, 1.
- Cao W, Wang D, Li J, et al., 2018. Brits: Bidirectional recurrent imputation for time series. *Advances in neural information processing systems*, 31.
- Che Z, Purushotham S, Cho K, et al., 2018. Recurrent neural networks for multivariate time series with missing values. *Scientific reports*, 8(1):6085.
- Choi J, Choi H, Hwang J, et al., 2022. Graph neural controlled differential equations for traffic forecasting. *Proceedings of the AAAI Conference on Artificial Intelligence*, 36(6):6367-6374.
- Chu C, Minami K, Fukumizu K, 2020. Smoothness and stability in gans. *arXiv preprint arXiv:200204185*, .
- Chung J, Gulcehre C, Cho K, et al., 2014. Empirical evaluation of gated recurrent neural networks on sequence modeling. *NIPS 2014 Workshop on Deep Learning*.
- Croitoru FA, Hondru V, Ionescu RT, et al., 2023. Diffusion models in vision: A survey. *IEEE Transactions on Pattern Analysis and Machine Intelligence*, to appear.
- Desai A, Freeman C, Wang Z, et al., 2021. Timevae: A variational auto-encoder for multivariate time series generation. *arXiv preprint arXiv:211108095*, .

- Dhariwal P, Nichol A, 2021. Diffusion models beat GANs on image synthesis. *Advances in Neural Information Processing Systems*, 34:8780-8794.
- Donahue C, McAuley J, Puckette M, 2018. Adversarial audio synthesis. *arXiv preprint arXiv:180204208*, .
- Esteban C, Hyland SL, Rättsch G, 2017. Real-valued (medical) time series generation with recurrent conditional gans. *arXiv preprint arXiv:170602633*, .
- Fortuin V, Baranchuk D, Rättsch G, et al., 2020. Gp-vae: Deep probabilistic time series imputation. *International conference on artificial intelligence and statistics*, p.1651-1661.
- Goel K, Gu A, Donahue C, et al., 2022. It’s raw! audio generation with state-space models. *International Conference on Machine Learning*, p.7616-7633.
- Gu A, Goel K, Re C, 2022. Efficiently modeling long sequences with structured state spaces. *International Conference on Learning Representations*.
- Harvey W, Naderiparizi S, Masrani V, et al., 2022. Flexible diffusion modeling of long videos. *Advances in Neural Information Processing Systems*, 35:27953-27965.
- Ho J, Jain A, Abbeel P, 2020. Denoising diffusion probabilistic models. *Advances in Neural Information Processing Systems*, 33:6840-6851.
- Ho J, Saharia C, Chan W, et al., 2022a. Cascaded diffusion models for high fidelity image generation. *Journal of Machine Learning Research*, 23(47):1-33.
- Ho J, Salimans T, Gritsenko A, et al., 2022b. Video diffusion models. *Advances in Neural Information Processing Systems*, 35:8633-8646.
- Hochreiter S, Schmidhuber J, 1997. Long short-term memory. *Neural computation*, 9(8):1735-1780.
- Hyvärinen A, Dayan P, 2005. Estimation of non-normalized statistical models by score matching. *Journal of Machine Learning Research*, 6(4).
- Kong Z, Ping W, Huang J, et al., 2021. Diffwave: A versatile diffusion model for audio synthesis. *International Conference on Learning Representations*, p.1-8.
- Li R, Li X, Gao S, et al., 2023. Graph convolution recurrent denoising diffusion model for multivariate probabilistic temporal forecasting. Working Paper, Unpublished.
- Li X, Thickstun J, Gulrajani I, et al., 2022. Diffusion-LM improves controllable text generation. *Advances in Neural Information Processing Systems*, 35:4328-4343.
- Li Y, Yu R, Shahabi C, et al., 2018. Diffusion convolutional recurrent neural network: Data-driven traffic forecasting. *International Conference on Learning Representations*, p.1-8.
- Li Y, Lu X, Wang Y, et al., 2022. Generative time series forecasting with diffusion, denoise, and disentanglement. *Advances in Neural Information Processing Systems*, 35:23009-23022.
- Li Y, Chen Z, Zha D, et al., 2021. Learning disentangled representations for time series. *arXiv preprint arXiv:210508179*, 1.

- Lim H, Kim M, Park S, et al., 2023. Regular time-series generation using sgm. *arXiv preprint arXiv:230108518*, 1.
- Liu M, Huang H, Feng H, et al., 2023. PriSTI: A conditional diffusion framework for spatiotemporal imputation. *International Conference on Data Engineering*, p.1-10.
- Luo C, 2022. Understanding diffusion models: A unified perspective. *arXiv preprint arXiv:220811970*, 1.
- Luo Y, Cai X, Zhang Y, et al., 2018. Multivariate time series imputation with generative adversarial networks. *Advances in neural information processing systems*, 31.
- Mogren O, 2016. C-rnn-gan: Continuous recurrent neural networks with adversarial training. *arXiv preprint arXiv:161109904*, .
- Mulyadi AW, Jun E, Suk HI, 2021. Uncertainty-aware variational-recurrent imputation network for clinical time series. *IEEE Transactions on Cybernetics*, 52(9):9684-9694.
- Nikolay S, Junyoung C, Mikolaj B, et al., 2022. Step-unrolled denoising autoencoders for text generation. *International Conference on Learning Representations*, p.1-8.
- van den Oord A, Dieleman S, Zen H, et al., 2016. Wavenet: A generative model for raw audio. *The 9th ISCA Speech Synthesis Workshop*, p.125.
- Osman MS, Abu-Mahfouz AM, Page PR, 2018. A survey on data imputation techniques: Water distribution system as a use case. *IEEE Access*, 6:63279-63291.
- Rasul K, Seward C, Schuster I, et al., 2021a. Autoregressive denoising diffusion models for multivariate probabilistic time series forecasting. *International Conference on Machine Learning*, p.8857-8868.
- Rasul K, Sheikh A, Schuster I, et al., 2021b. Multi-variate probabilistic time series forecasting via conditioned normalizing flows. *International Conference on Learning Representations*, p.1-8.
- Ronneberger O, Fischer P, Brox T, 2015. U-net: Convolutional networks for biomedical image segmentation. *Medical Image Computing and Computer-Assisted Intervention*, p.234-241.
- Saremi S, Hyvarinen A, 2019. Neural empirical Bayes. *Journal of Machine Learning Research*, 20(181):1-23.
- Simeunović J, Schubnel B, Alet PJ, et al., 2021. Spatio-temporal graph neural networks for multi-site pv power forecasting. *IEEE Transactions on Sustainable Energy*, 13(2):1210-1220.
- Sohl-Dickstein J, Weiss E, Maheswaranathan N, et al., 2015. Deep unsupervised learning using nonequilibrium thermodynamics. *International Conference on Machine Learning*, p.2256-2265.
- Song Y, Ermon S, 2019. Generative modeling by estimating gradients of the data distribution. *Advances in Neural Information Processing Systems*, 32:11918-11930.
- Song Y, Garg S, Shi J, et al., 2020. Sliced score matching: A scalable approach to density and score estimation. *Uncertainty in Artificial Intelligence*, p.574-584.

- Song Y, Sohl-Dickstein J, Kingma DP, et al., 2021. Score-based generative modeling through stochastic differential equations. *International Conference on Learning Representations*, p.1-8.
- Tashiro Y, Song J, Song Y, et al., 2021. CSDI: Conditional score-based diffusion models for probabilistic time series imputation. *Advances in Neural Information Processing Systems*, 34:24804-24816.
- Vahdat A, Kautz J, 2020. NVAE: A deep hierarchical variational autoencoder. *Advances in Neural Information Processing Systems*, 33:19667-19679.
- Vaswani A, Shazeer N, Parmar N, et al., 2017. Attention is all you need. *Advances in Neural Information Processing Systems*, 30:6000-6010.
- Vincent P, 2011. A connection between score matching and denoising autoencoders. *Neural Computation*, 23(7):1661-1674.
https://doi.org/10.1162/NECO_a_00142
- Wen H, Lin Y, Xia Y, et al., 2023. Diffstg: Probabilistic spatio-temporal graph forecasting with denoising diffusion models. *arXiv preprint arXiv:230113629*, 1.
- Wu H, Xu J, Wang J, et al., 2021. Autoformer: Decomposition transformers with auto-correlation for long-term series forecasting. *Advances in Neural Information Processing Systems*, 34:22419-22430.
- Xiao Z, Kreis K, Vahdat A, 2021. Tackling the generative learning trilemma with denoising diffusion gans. *arXiv preprint arXiv:211207804*, .
- Xu T, Wenliang LK, Munn M, et al., 2020. Cot-gan: Generating sequential data via causal optimal transport. *Advances in neural information processing systems*, 33:8798-8809.
- Yan T, Zhang H, Zhou T, et al., 2021. Scoregrad: Multivariate probabilistic time series forecasting with continuous energy-based generative models. *arXiv preprint arXiv:210610121*, 1.
- Yang L, Zhang Z, Song Y, et al., 2022a. Diffusion models: A comprehensive survey of methods and applications. *arXiv preprint arXiv:220900796*, 1.
- Yang R, Srivastava P, Mandt S, 2022b. Diffusion probabilistic modeling for video generation. *arXiv preprint arXiv:220309481*, 1.
- Yi X, Zheng Y, Zhang J, et al., 2016. St-mvl: filling missing values in geo-sensory time series data. *Proceedings of the 25th International Joint Conference on Artificial Intelligence*.
- Yoon J, Zame WR, van der Schaar M, 2018. Estimating missing data in temporal data streams using multi-directional recurrent neural networks. *IEEE Transactions on Biomedical Engineering*, 66(5):1477-1490.
- Yu P, Xie S, Ma X, et al., 2022. Latent diffusion energy-based model for interpretable text modelling. *International Conference on Machine Learning*, 162:25702-25720.

Zhou X, Shen H, Ye J, 2011. Integrating outlier filtering in large margin training. *Journal of Zhejiang University Science C*, 12(5):362-370.
<https://doi.org/10.1631/jzus.C1000361>