

Probabilistic Decomposition Transformer for Time Series Forecasting

Junlong Tong, Liping Xie*, Wankou Yang, Kanjian Zhang

School of Automation, Southeast University

Abstract

Time series forecasting is crucial for many fields, such as disaster warning, weather prediction, and energy consumption. The Transformer-based models are considered to have revolutionized the field of sequence modeling. However, the complex temporal patterns of the time series hinder the model from mining reliable temporal dependencies. Furthermore, the autoregressive form of the Transformer introduces cumulative errors in the inference step. In this paper, we propose the probabilistic decomposition Transformer model that combines the Transformer with a conditional generative model, which provides hierarchical and interpretable probabilistic forecasts for intricate time series. The Transformer is employed to learn temporal patterns and implement primary probabilistic forecasts, while the conditional generative model is used to achieve non-autoregressive hierarchical probabilistic forecasts by introducing latent space feature representations. In addition, the conditional generative model reconstructs typical features of the series, such as seasonality and trend terms, from probability distributions in the latent space to enable complex pattern separation and provide interpretable forecasts. Extensive experiments on several datasets demonstrate the effectiveness and robustness of the proposed model, indicating that it compares favorably with the state of the art.

Keywords: Time series forecasting, Variational inference, Series decomposition

1. Introduction

In the era of big data, the explosive growth of data has forced the development of data mining. Time series information, as one of the most common data, has captured widespread attention for the mining of its intrinsic patterns. Time series data mining refers to discovering potential features to provide decision making, such as prediction, classification, and anomaly detection, by learning historical data in chronological order. Among them, the importance and application of time series forecasting have been demonstrated in modern society such as commodity demand forecasting [1], energy consumption [2], traffic planning [3], and financial analysis [4].

Classical statistical methods have been well used for time series forecasting, such as ARIMA [5] models and state space models (SSMs) [6]. These methods usually incorporate prior knowledge of the time series, such as trend and seasonality characteristics, and use them for decision making and interpretation. However, statistical-based methods fail to capture the relationship between the covariates and the target series, limiting the effectiveness of predicting intricate time series.

In recent years, deep prediction models have been significantly developed to tackle large-scale complex time series forecasting problems. For example, models based on recurrent neural networks (RNNs) [7, 8] and self-attention mechanisms

[9, 10] are widely used for time series forecasting. Models based on the self-attention mechanism simultaneously calculate the attention scores between any two time points, assigning different weights according to the importance of each part of the input data. Compared with the circular feature of RNN, the self-attention mechanism-based model processes the historical time series data at once and realizes the information interaction at different moments, which can maintain the long-term dependencies. However, forecasting models based on RNN and self-attention mechanism follow an autoregressive form and usually employ the teacher-forcing strategy [11, 12], i.e., providing ground truth at each moment in the training phase, to improve convergence and generalization. However, the strategy causes inconsistencies between the training and inference phases, leading to exposure bias [13] in the inference phase.

In addition, the complex temporal patterns of the time series prevent the models from mining reliable temporal dependencies. To separate the complex patterns, the concept of series decomposition [14, 15] is introduced into time series analysis, which assumes that the sequence consists of components based on prior knowledge and seeks to separate them. However, due to the specificity of the forecasting task, only the input historical series can be pre-processed, ignoring the interaction with the forecast information and lacking flexibility.

In this work, we propose to combine the strengths of the Transformer architecture and conditional generative model, for hierarchical and interpretable probabilistic forecasting. The Transformer [9] is used for temporal feature extraction and primary probabilistic forecasting, where the probability distribu-

*Corresponding author

Email addresses: jltong@seu.edu.cn (Junlong Tong),
lpxie@seu.edu.cn (Liping Xie), wkyang@seu.edu.cn (Wankou Yang),
kjzhang@seu.edu.cn (Kanjian Zhang)

tion parameters of the time series are forecast by an autoregressive process. In addition, the probability distribution parameters are used as conditional information for probabilistic encoding and reconstruction of the prediction using a variational inference [16, 17] generative model. The hierarchical probabilistic forecasting method can effectively mitigate the exposure bias brought by the autoregressive Transformer model by introducing a conditional generative model to constrain the primary forecasting results at the sequence level. In addition, the probabilistic decoder is designed based on the prior knowledge of the time series to reconstruct the subsequence of different features from the latent space to achieve the separation of intricate patterns and interpretable forecasts, where the probabilistic decoder combines historical sequences and primary forecasting result for information interaction.

The principal contributions of this work can be summarized as follows:

- An efficient time series forecasting model called Probabilistic Decomposition Transformer is proposed in this work, where the model combines the Transformer architecture and generative model based on variational inference for hierarchical probabilistic forecasting.
- We design a novel probabilistic decomposition framework, where the probabilistic decoder reconstructs typical features of sequences from probability distributions in the latent space to achieve separation of intricate temporal patterns and provide interpretable forecasts.
- The experiments demonstrate that the proposed model is effective in reducing the exposure bias of autoregressive forecasting, showing that it compares favorably with the state-of-the-art models.

2. Related Work

2.1. Time Series Forecasting

The early methods used for time series forecasting are mainly statistical models, such as ARIMA [5] and exponential smoothing [18]. However, statistical models have difficulties capturing complex temporal patterns, which limits their application to complex scenarios. Recently, deep neural networks have been applied to time series forecasting thanks to the powerful feature representation capabilities and scalable structure. The most prominent of the deep models is the recurrent neural networks (RNNs) [19, 20]. DeepAR [21] presented an autoregressive RNNs method for probabilistic distribution modeling. DSSM [22] combined RNNs with a state space model for probabilistic time series forecasting. Besides, convolutional neural networks (CNNs) also have been proposed for time series forecasting. LSTNet [23] combined CNNs with skip-RNN to extract local patterns among variables and capture long-term dependency based on time series trends. Bai et al. [24] proposed the temporal convolution network (TCN), a sequence model based on causal convolutions and dilated convolutions, which has been used for time series forecasting [25, 26].

In recent years, Transformer-based [9, 27] models have attracted the attention of researchers. Most of the work is dedicated to modifying the self-attention mechanism [10, 28, 29] or to improving the Transformer architecture [30, 31]. For example, Li et al. [10] first applied Transformer to time series forecasting, where a convolution Transformer is proposed to improve local processing power and a log-sparse strategy is designed for breaking the memory bottleneck. SSDNet [32] conducted the Transformer to learn the parameters of the state space model to provide probabilistic and interpretable forecasts. In this work, Transformer is employed to extract temporal patterns and implement primary probabilistic forecasting.

2.2. Decomposition of Time Series

Time series typically contain multiple underlying patterns that increase the complexity of prediction. Time series decomposition [33, 15] is regarded as an effective method for pattern separation, which is to decompose a series into several components, each of which corresponds to a particular pattern. For time series forecasting tasks, a variety of decomposition strategies have been proposed for mining historical volatility over time. For instance, Prophet [34] provides a trend-seasonality decomposition, and DeepGLO[25] implements a global matrix factorization model for high-dimensional time series.

In recent years, learning-based decomposition strategies have attracted the attention of researchers. For instance, Oreshkin et al. [35] proposed the N-Beats model with basis expansion, an interpretable prediction framework with fully-connected layer structures. Nguyen et al. [36] proposed a temporal latent autoencoder for the nonlinear decomposition of multivariate time series. Furthermore, in Autofomer[28] and FEDformer[29], decomposition based on average pooling operation is conducted to highlight the trend term, where the decomposition module is considered as the inner block of forecasting model to decompose series progressively. This paper proposes a new decomposition idea, which obtains primary prediction results through Transformer model, then uses conditional generation model for probabilistic decomposition to reconstruct subsequences of different features from potential space. The process effectively realizes the interaction between historical data and prediction data, breaks the information barrier, and improves the reliability of decomposition results.

2.3. Generative Model for Time Series

Generative models have been extensively used for time series data mining tasks, such as time series imputation [37, 38], time series data generation [39, 40, 41]. For forecasting tasks, there have been many generative models based on variational inference [16, 42] introduced to conditional probability modeling. Fortuin et al. [43] proposed a VAE-based method for deep probability time series imputation. Tang et al. [44] proposed to combine the SSM model with the attention mechanism to achieve non-autoregressive probability forecasting by conditional variational inference. Li et al. [46] proposed a reformulated VAE framework for time series disease forecasts. Besides, a recent approach based on diffusion model [45] is proposed to

autoregressive time series forecasting. In this work, a conditional generative model is employed to tackle the conditional probability of the primary forecasting result.

3. Preliminaries

3.1. Problem Definition

Let $\{Y_{i,1:t_0}\}_{i=1}^N$ represent a set of N related univariate time series, where $Y_{i,1:t_0} = [y_{i,1}, y_{i,2}, \dots, y_{i,t_0}] \in \mathbb{R}^{1 \times t_0}$ and $y_{i,t}$ is the value of i th time series at time t . In addition, assume that $X_{i,1:t_0} \in \mathbb{R}^{k \times t_0}$ denotes the k -dimensional covariate series corresponding to $Y_{i,1:t_0}$, where the covariates can be static (i.e., time-independent features such as serial number) or dynamic (i.e., time-related features such as time period of the day).

Specifically, the task of this work is to model the conditional distribution $p(Y_{i,t_0+1:t_0+\tau} | Y_{i,1:t_0}, X_{i,1:t_0+\tau}; \Phi)$, where the future covariates are known. The conditional distribution can be reduced by a step-by-step forecasting task, where the problem can be denoted as follows:

$$\begin{aligned} p(Y_{i,t_0+1:t_0+\tau} | Y_{i,1:t_0}, X_{i,1:t_0+\tau}; \Phi) \\ = \prod_{t=t_0+1}^{t_0+\tau} p(Y_{i,t} | Y_{i,1:t-1}, X_{i,1:t}; \Phi), \end{aligned} \quad (1)$$

where Φ is parameters of the model trained on all N related univariate time series, and the input of the proposed model at time t is a concatenated vector of Y_{t-1} and X_t . Besides, the time range $[1, t_0]$ and $[t_0 + 1, t_0 + \tau]$ are referred as conditional range and prediction range respectively, to be consistent with the literature [10, 21], where t_0 represents the forecasting start moment and τ represents the forecasting horizon. Moreover, owing to all univariate time series being trained with the same model, the identification i to distinguish different time series will be omitted in the latter.

3.2. The Transformer

Structurally, Transformer-based models [9, 47] are encoder-decoder architectures based mainly on the self-attention mechanism, which enables the models to capture long-term dependencies and focus on important patterns. In practice, multi-headed self-attention is usually applied to the Transformer to improve the model fitting performance. For the self-attention model with H head, the outputs O_1, \dots, O_H are concatenated and linearly projected to O , i.e., $O = \text{concat}(O_1, \dots, O_H)W^O$, where W^O is the parameter of the linear projection. For $h \in [1, H]$, multi-head self-attention output O_h can be expressed as follows:

$$\begin{aligned} O_h &= \text{Attention}(Q_h, K_h, V_h) \\ &= \text{softmax}\left(\frac{Q_h K_h^T}{\sqrt{d_k}}\right)V_h, \end{aligned} \quad (2)$$

where $Q_h = QW_h^O$, $K_h = KW_h^K$, and $V_h = VW_h^V$ are the projections of query, key, and value with learnable parameters W_h^O , W_h^K , and W_h^V , respectively.

Compared with the recurrent neural network [48, 49] that passes memory states step by step, the self-attention mechanism calculates the attention scores of any two moments simultaneously, which improves the learning ability of long-term dependencies as well as optimizes parallel performance. Besides, a masking mechanism is introduced to the decoder to obscure future information for preventing information leakage.

3.3. Variational Autoencoder

Variational inference is a parameter estimation method for dealing with the conditional probability containing hidden variables by introducing a variational distribution, usually implemented by the variational autoencoder (VAE). The VAE [16] is a deep generative model designed to learn the probability distribution of existing datasets and generate new data through sampling from the learned data distribution. Since the form of the probability distribution of existing data is unknown, VAE encodes the input data into the latent space, constructs the probability distribution of the latent variables, samples from the latent variables and reconstructs new data.

4. Probabilistic Decomposition Transformer

4.1. Model Architecture

In brief, the aim of the Probabilistic Decomposition Transformer (PDTrans) is to build a hierarchical probabilistic forecasting model to provide accurate probabilistic and interpretable forecasts.

The architecture of the proposed model is illustrated in Figure 1, which contains two important components: Transformer and conditional generative model. The Transformer learns temporal patterns and implements autoregressive probabilistic forecasting, a process that outputs the parameters of the probability distribution to provide primary forecasting. The function of the conditional generative model is twofold. The first function is to model the primary forecasting distribution through variational inference to achieve hierarchical forecasting, which can mitigate the impact of exposure bias in the autoregressive prediction process. Reconstructing typical features of sequence from the probability distribution of the latent space is another function, where a probabilistic decoder is employed to achieve the separation of complex patterns and provide interpretable forecasts.

4.2. Autoregressive Probabilistic Forecasts

We attempt to model the conditional probabilities shown in Equation 1 and assume that each conditional probability obeys a specific probability distribution that can be represented by a learnable likelihood function, as shown in the following:

$$p(Y_t | Y_{1:t-1}, X_{1:t}; \Phi) = l(Y_t | \Phi_t), \quad (3)$$

where $l(Y_t | \Phi_t)$ is the likelihood and Φ_t is the parameter of the likelihood at moment t . The likelihood should be selected based on the statistical properties of the data to accurately approximate the true distribution. For real-world datasets, the Gaussian likelihood is usually chosen to approximate the conditional probability distribution.

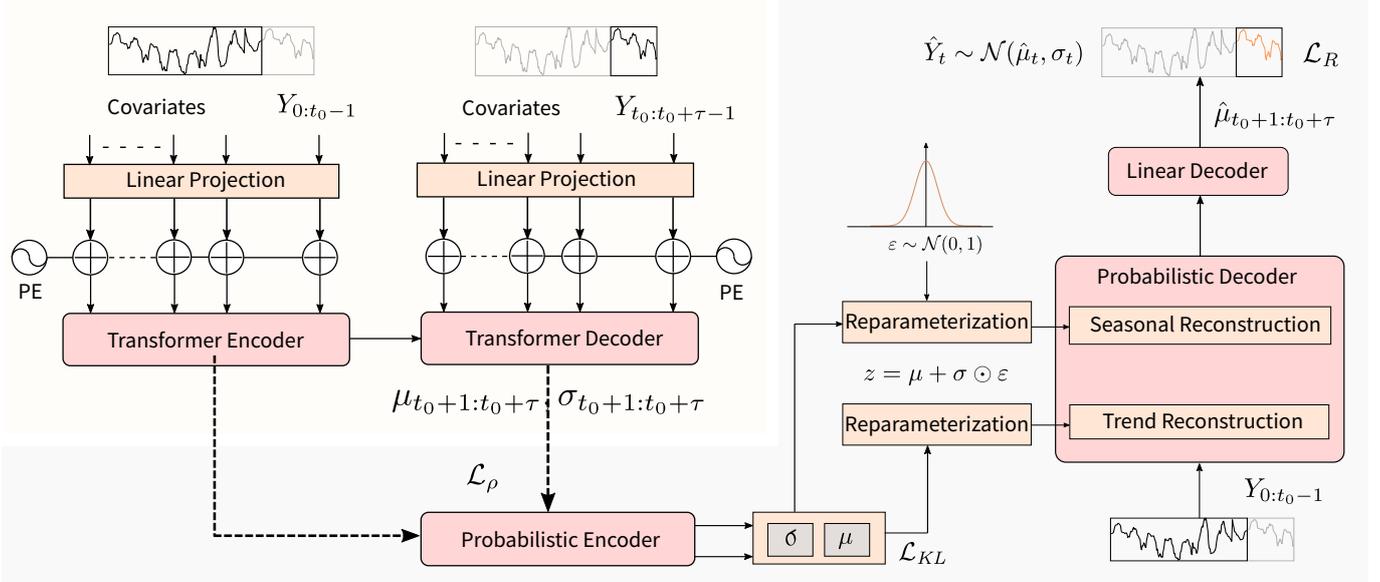


Figure 1: Architecture of the PDTrans model. (left) Transformer model for primary autoregressive probabilistic forecasts. (right) Conditional generative model for sequence-level probabilistic forecasts and pattern separation.

Since the Transformer has made a splash in the field of sequence modeling, this work employs Transformer to tackle the autoregressive probabilistic forecasting problem. The encoder of Transformer takes the linear projection of history sequence and covariates as input, where the linear projection is applied to learnable embedding, as shown in Figure 1 (left). Besides, the position information is represented by adding the position encode (PE) to the output of the embedding layer. The decoder of Transformer outputs the parameters of the likelihood function at each moment in an autoregressive way, which represents the approximation to the conditional probability distribution.

In practice, MLP is introduced to perform an affine transformation on the decoder output to constrain the range of parameter values. For example, for Gaussian likelihood, the parameters can be expressed as $\Phi_t = (\mu_t, \sigma_t)$, where the μ_t and σ_t represent the mean and standard deviation of the likelihood function, respectively. Then the likelihood and the parameters can be shown in Equation 4 and 5, respectively.

$$l(Y_t|\mu_t, \sigma_t) = \frac{\exp(-(Y_t - \mu_t)^2/2\sigma_t^2)}{\sqrt{2\pi\sigma_t^2}}, \quad (4)$$

$$\begin{aligned} \mu_t &= W_\mu^T f_t + b_\mu, \\ \sigma_t &= \text{softplus}(w_\sigma^T f_t + b_\sigma), \end{aligned} \quad (5)$$

where f_t is the output of Transformer at moment t . The softplus function is chosen to ensure that the model generates positive standard deviations.

In the training phase, the historical and predicted sequences are simultaneously input to the model, where the parallel forecasting is achieved through a masking mechanism. Owing to the lack of ground truth during inference, the forecasted values need to be fed back to the Transformer to achieve sequence

forecasting. We sample \hat{Y}_t as the ground truth according to $\hat{Y}_t \sim l(Y_t|\mu_t, \sigma_t)$ and feed the samples to the model.

4.3. Conditional Generative Forecasts

In order to mitigate the impact of exposure bias and provide the separation of intricate patterns, we introduce a conditional generative model to tackle the primary forecasting result of the Transformer, which provides hierarchical forecasting and interpretable forecasting.

Hierarchical forecasting represents that the conditional generative model performs non-autoregressive forecasting of the Transformer prediction results, which brings the forecasted distribution closer to the true distribution of the model. Pattern separation is implemented by the sequence decomposition module of the probabilistic decoder, which reconstructs trend and seasonality terms from latent space. In addition, the primary forecasting results provided by the Transformer enable the interaction of historical sequences and future information, enhancing the information extraction capability of the decomposition module.

Specifically, we are interested in the probabilistic model parameterized by θ of the form:

$$\begin{aligned} & p_\theta(Y_{t_0+1:t_0+\tau}|Y_{1:t_0}) \\ &= \int_z p_\theta(Y_{t_0+1:t_0+\tau}|z) p_\theta(z|Y_{1:t_0}) dz, \end{aligned} \quad (6)$$

where $Y_{t_0+1:t_0+\tau}$ represents the target series from moment t_0+1 to $t_0+\tau$, z represents the latent variables, and $p_\theta(z|Y_{1:t_0})$ represents the prior distribution.

The marginal likelihood and posterior density are intractable that require approximate posterior inference necessarily. We follow the framework of stochastic variational inference [16, 42] and suppose that the variational posterior is parameterized

by ϕ . Then the evidence lower bound (ELBO) of the model can be written as follows:

$$\begin{aligned} & \log p_\theta(Y_{t_0+1:t_0+\tau}|Y_{1:t_0}) \\ & \geq -D_{KL}(q_\phi(z|Y_{1:t_0}, Y_{t_0+1:t_0+\tau})||p_\theta(z|Y_{1:t_0})) + \\ & \quad \mathbb{E}_{q_\phi(z|Y_{1:t_0}, Y_{t_0+1:t_0+\tau})}[\log p_\theta(Y_{t_0+1:t_0+\tau}|z, Y_{1:t_0})], \end{aligned} \quad (7)$$

where q_ϕ is an approximation of the true posterior and D_{KL} represents the Kullback-Leibler (KL) divergence. The derivation for the ELBO can be found in Appendix. Besides, we choose a factor β for the KL term to balance capacity of latent variables and independence constraints with reconstruction accuracy, following relevant work [50].

However, the target series $Y_{t_0+1:t_0+\tau}$ is unknown under the forecasting scenario, which hinders the optimization of the above problem. Therefore, the generative model in this work actually optimizes the likelihood function output by the Transformer and reconstructs the probability distribution that is close enough to the true distribution. Then the ELBO can be denoted as follows:

$$\begin{aligned} ELBO = & -D_{KL}(q_\phi(z|Y_{1:t_0}, \mu_{t_0+1:t_0+\tau})||p_\theta(z|Y_{1:t_0})) \\ & + \mathbb{E}_{q_\phi(z|Y_{1:t_0}, \mu_{t_0+1:t_0+\tau})}[\log p_\theta(\mu_{t_0+1:t_0+\tau}|z, Y_{1:t_0})], \end{aligned} \quad (8)$$

where the $\mu_{t_0+1:t_0+\tau}$ represents the parameters of the likelihood function from moment $t_0 + 1$ to $t_0 + \tau$.

4.3.1. Inference Model

For the generative forecast model, the likelihood parameters are first mapped by the inference model into the latent variable space, where the input sequence can be represented by the latent variables.

The inference model (also known as recognition model) $q_\phi(z|Y_{1:t_0}, \mu_{t_0+1:t_0+\tau})$ is a conditional Bayesian network, which is adopted to approximate the intractable posterior distribution $p_\theta(z|Y_{1:t_0})$. Besides, the KL divergence is applied to assess the similarity between the inference model and the true posterior. As shown in the right of Figure 1, the inferential model can be implemented by the probabilistic encoder. We utilize MLP to achieve the probabilistic encoder, where the input of the probabilistic encoder consists of the historical series $Y_{1:t_0}$ and the parameters of the likelihood function $\mu_{t_0+1:t_0+\tau}$.

The probabilistic encoder outputs the parameters of a Gaussian distribution and samples the latent variable z by a reparameterization trick [51]. The reparameterization trick converts the representation of the random variable z into a deterministic part and a stochastic part, shown as follows:

$$\begin{aligned} z &= \mu + \sigma \odot \varepsilon, \\ \varepsilon &\sim \mathcal{N}(0, 1), \end{aligned} \quad (9)$$

where μ and σ are parameters of Gaussian distribution generated by the probabilistic encoder, ε is random variable sampled by standard Gaussian distribution.

4.3.2. Generative Model

The generative model refers to the generation of forecasting values according to observed variables and latent variables, which can be expressed as $p_\theta(\mu_{t_0+1:t_0+\tau}|z, Y_{1:t_0})$.

In this work, the generator is employed to reconstruct the likelihood parameters, and we want the reconstructed likelihood function to be close enough to the true conditional distribution. Different from the reconstructed input of the original VAE, the model in this paper aims at approximating the ground truth. Therefore, we choose the negative log-likelihood function as the reconstruction loss in the generative model.

In addition, another task of the generative model is pattern separation for interpretable prediction. In the context of forecasts, the probabilistic encoder maps the primary forecasts result to the latent space, which can compensate for the lack of future information and improve the performance of the model to decompose the target series. The probability decoder reconstructs typical features of the sequence, such as seasonality and trend terms, from Gaussian distributions in the latent space to achieve pattern separation and provide interpretable forecasts, where the historical sequence is fed into the decoder as conditional information. As shown in the right of Figure 1.

Let $\hat{\mu}_t$ denotes the reconstructed likelihood parameters corresponding to the target time series at time t , where $Y_t \sim \mathcal{N}(\hat{\mu}_t, \sigma_t^2)$. We introduce the probabilistic decoder to decompose the above distribution into two Gaussian distributions with the following characteristics:

$$\begin{aligned} Y_t^{trend} &\sim \mathcal{N}(\mu_t^{trend}, \sigma_t^2/2), \\ Y_t^{seasonal} &\sim \mathcal{N}(\mu_t^{seasonal}, \sigma_t^2/2), \\ \hat{\mu}_t &= \mu_t^{trend} + \mu_t^{seasonal}, \end{aligned} \quad (10)$$

where μ_t^{trend} and $\mu_t^{seasonal}$ represent likelihood parameters of the trend term and seasonality term at time t , respectively.

Specifically, the probability decoder takes the latent variable as input and obtains the decomposition features through the trend feature extraction and seasonality feature extraction modules. The trend term can be extracted explicitly by convolution operation that can simulate the moving average to smooth out periodic fluctuations and highlight the trends. The prediction results are obtained by linear decoding of trend term and seasonality term, where the seasonality term can be extracted implicitly by MLP. Let AvgPooling represents the average pooling operation for moving average and LinearDecoder denotes a linear decoder, then the process is shown as follows:

$$\begin{aligned} \mu_{t_0+1:t_0+\tau}^{trend} &= \text{AvgPooling}(\text{MLP}(z)), \\ \mu_{t_0+1:t_0+\tau}^{seasonal} &= \text{MLP}(z), \\ \hat{\mu}_{t_0+1:t_0+\tau} &= \text{LinearDecoder}(\mu_{t_0+1:t_0+\tau}^{trend}, \mu_{t_0+1:t_0+\tau}^{seasonal}). \end{aligned} \quad (11)$$

Then, given likelihood parameters $\hat{\mu}_t$ and σ_t , the forecasting result can be obtained by sampling the Gaussian distribution $\hat{Y}_t \sim \mathcal{N}(\hat{\mu}_t, \sigma_t^2)$.

4.4. Joint Learning

We define several loss terms to train the proposed model jointly, where the Transformer contains the negative log-likelihood (NLL) loss and the conditional generative model contains the KL divergence and reconstruction loss.

4.4.1. Negative Log-Likelihood loss

The Transformer predicts the conditional distributions by maximum likelihood estimation. The negative log-likelihood is chosen as the loss function, where the minimization loss function is equivalent to the maximum likelihood estimation. For Gaussian likelihood, the loss of NLL is shown as follows:

$$\mathcal{L}_{NLL} = \sum_t l(Y_t | \mu_t, \sigma_t) = \sum_t \frac{(Y_t - \mu_t)^2}{2\sigma_t^2} + \log \sigma_t + \text{Const}, \quad (12)$$

where the detail can be found in [21].

4.4.2. KL loss

As in VAE, we regularize the latent space by encouraging it to be similar to a standard Gaussian distribution with μ the null vector and σ the identity matrix. We minimize the KL divergence between the encoder distribution and the true posterior.

$$\begin{aligned} \mathcal{L}_{KL} &= -D_{KL}(q_\phi(z|Y_{1:t_0}, \mu_{t_0+1:t_0+\tau}) || p_\theta(z|Y_{1:t_0})) \\ &= \frac{1}{2} \sum_{j=1}^N (1 + \log \sigma_j^2 - \mu_j^2 - \sigma_j^2), \end{aligned} \quad (13)$$

where N is the number of the latent variables, μ_j and σ_j are element of μ and σ , respectively.

4.4.3. Reconstruction loss

One of the goals of the generative model is to maximize the expectation term in evidence lower bound, i.e., maximum $\mathbb{E}_{q_\phi(z|Y_{1:t_0}, \mu_{t_0+1:t_0+\tau})} [\log p_\theta(\mu_{t_0+1:t_0+\tau} | z, Y_{1:t_0})]$. The maximum problem is equivalent to minimizing the reconstruction loss. For the probabilistic decoder, the reconstruction loss represents the error of the generative model output with respect to the likelihood function, which can be denoted as $\mathcal{L}_R(\hat{\mu}_t, \mu_t)$. Considering that the likelihood function obtained from the Transformer is used to approximate the true distribution, we choose to directly measure the loss between the likelihood of the reconstruction and the true distribution, which can be denoted as $\mathcal{L}_R(\hat{\mu}_t, Y_t)$. In experiments, $\mathcal{L}_R(\hat{\mu}_t, Y_t)$ can be chosen as the NLL, the above scheme was observed to achieve better forecasting performance.

The total loss of the PDTrans is defined as the weighted summation of different terms: $\mathcal{L} = \gamma \mathcal{L}_{NLL} + \beta \mathcal{L}_{KL} + \mathcal{L}_R$, where β and γ are trade-off coefficients.

5. Experiments

5.1. Datasets

In this work, five public datasets are used for model performance evaluation, i.e., Electricity, Traffic, Solar, Exchange, and M4-Hourly. The general statistics are listed in Table 1.

The Electricity dataset consists of the electricity consumption of 370 users at a 15-minute resolution from 01/01/2011 to 07/09/2014, where in this work all series are aggregated into hourly intervals. The Traffic dataset is composed of the road

Table 1: General information about datasets, where the 3rd column represents the sequence length of each variable.

Datasets	Resolution	Length	Variables
Electricity	1 hour	32,304	370
Traffic	1 hour	4,049	963
Solar	1 hour	4,832	137
Exchange	1 day	7,588	8
M4-Hourly	1 hour	748/1,008	414

Table 2: Hyperparameters for PDTrans, where N is number of layers, h is number of head, d_{model} is dimensionality of model, d_{ff} is dimensionality of inner-layer, λ and β represent trade-off coefficients, k is kernel size of 1-D convolutions.

	N	h	d_{model}	d_{ff}	λ	β	k
Electricity	3	8	160	2048	1	1	5
Traffic	3	8	160	2048	1	1	3
Solar	3	8	16	640	1	1	3
Exchange	3	4	160	640	1	1	5
M4-Hourly	3	8	160	2048	1	1	3

occupancy rate, between 0 and 1, where the dataset contains 963 related variables with hourly resolution from 02/01/2008 to 30/03/2009. The Solar dataset consists of hourly interval solar power production data collected from 137 photovoltaic plants in Alabama from January to August 2006. The Exchange dataset contains daily exchange rate records from January 1990 to December 2016 for 8 countries. The M4-Hourly dataset includes 414 time series of hourly intervals from M4 competition [52], where the training and test set have been provided.

For covariates, we consider both static features and dynamic features, where the static features refer to the serial ID and the dynamic features refer to the time-related features such as time period and relative positions. The time period contains the encoding of hour and day for hourly interval datasets, and the encoding of day and month for daily interval datasets. The relative positions can be described as the distance to the first value of the series for all datasets.

5.2. Evaluation Metrics And Experimental Setup

Following [21, 10], we evaluate the model using ρ -quantile loss with $\rho \in (0, 1)$, which is shown as follows:

$$Q_\rho(Y, \hat{Y}) = \frac{2 \sum_t P_\rho(y_t, \hat{y}_t)}{\sum_t |y_t|}, \quad (14)$$

$$P_\rho(y_t, \hat{y}_t) = (y_t - \hat{y}_t)(\rho I_{\hat{y}_t > y_t} - (1 - \rho) I_{\hat{y}_t \leq y_t}), \quad (15)$$

where y_t is the ground truth, \hat{y}_t represents the predicted distribution with ρ -quantile, and I is a boolean function. We report $\rho_{0.5}$ and $\rho_{0.9}$ metrics to be consistent with the literature [21, 10].

We utilize Vaswani Transformer [9] for primary forecasts, where the number of layers of both encoder and decoder is set to 3, and the head of the attention mechanism is set to 4 and 8 for

Table 3: $\rho_{0.5}/\rho_{0.9}$ metrics for the short-term (1d ahead) and long-term (7d ahead) forecasting scenarios on Electricity and Traffic datasets for autoregressive probabilistic forecasting models. \diamond denotes results from [32].

	DeepAR	DeepSSM	ConvTrans	SSDNet	PDTrans(our)
Elect _{1d}	0.075 \diamond / 0.040 \diamond	0.083 \diamond / 0.056 \diamond	0.059 \diamond / 0.034 \diamond	0.068 \diamond / 0.033 \diamond	0.058 / 0.030
Elect _{7d}	0.082 \diamond / 0.053 \diamond	0.085 \diamond / 0.052 \diamond	0.070 \diamond / 0.044 \diamond	0.079 / 0.042	0.068 / 0.038
Traffic _{1d}	0.161 \diamond / 0.099 \diamond	0.167 \diamond / 0.113 \diamond	0.122 \diamond / 0.081\diamond	0.153 / 0.101	0.113 / 0.088
Traffic _{7d}	0.179 \diamond / 0.105 \diamond	0.168 \diamond / 0.114 \diamond	0.139 \diamond / 0.094 \diamond	0.161 / 0.109	0.126 / 0.093

Table 4: $\rho_{0.5}/\rho_{0.9}$ metrics of various methods on five real datasets. \diamond indicates results reported by [32].

	Electricity	Traffic	Solar	Exchange	M4-Hourly
Prophet	0.112 / 0.055	0.183 / 0.137	0.268 / 0.169	0.017 / 0.013	0.102 / 0.038
DeepAR	0.075 \diamond / 0.040 \diamond	0.161 / 0.099	0.222 \diamond / 0.093 \diamond	0.014 \diamond / 0.009 \diamond	0.090 / 0.030
DeepSSM	0.083 \diamond / 0.056 \diamond	0.167 / 0.133	0.223 \diamond / 0.181 \diamond	0.014 \diamond / 0.012 \diamond	0.044 / 0.026
ConvTrans	0.059 \diamond / 0.034 \diamond	0.122 \diamond / 0.081\diamond	0.210 \diamond / 0.082 \diamond	0.017 / 0.008	0.067 / 0.025
N-Beats	0.061 \diamond / -	0.114 / -	0.212 / -	0.018 / -	0.025 / -
Informer	0.068 \diamond / -	0.122 / -	0.215 \diamond / -	0.014 \diamond / -	- / -
Autoformer	0.083 / -	0.120 / -	0.211 / -	0.013 / -	- / -
SSDNet	0.068 \diamond / 0.033 \diamond	0.166 / 0.106	0.209 \diamond / 0.074 \diamond	0.013 \diamond / 0.006\diamond	0.038 / 0.023
PDTrans(our)	0.058 / 0.030	0.113 / 0.088	0.205 / 0.073	0.011 / 0.006	0.033 / 0.023

different datasets. The dimensionality of the inner layer d_{ff} is set to 2048 for Electricity, Traffic, and M4-Hourly datasets, and set to 640 for Solar and Exchange datasets. The dimensionality of the model d_{model} is set to 16 for the Solar dataset and set to 160 for the others. For the primary forecasts model, we utilize learnable position and covariates embedding. The embedding dimensionality is set to 20 and 30. For the generative model, both of the trade-off coefficients β and γ are set to 1 for all datasets. The dimensionality of latent space is set to 20 for all datasets. The kernel size of the 1-D convolutions in the trend extraction module is set to 3 and 5 for different datasets. The main hyperparameters used in this work are listed in Table 2.

Furthermore, the split of the datasets is followed by the literature [21, 10, 32]. We utilize one week of data from 9/1/2014 on Electricity and 6/16/2008 on Traffic as test sets. For the Solar dataset, the last 7 days in August are used as test set. In addition, we leave one week of data for the validation set and the rest of the data as the training set for the above datasets. For the Exchange dataset, only the weekday data is considered, where 90% of the data before February 2, 2015 is used as the training set and 10% of the data is used as the validation set. The rest of the data after February 2, 2015 is used as the test set. The M4-Hourly dataset has been provided with a division of the training, validation, and test sets.

For all datasets, the Adam [53] optimizer is employed to optimize the model with a learning rate of 0.001, where the learning rate is set to decrease by 20% every two epochs. Besides, we set the maximum number of epochs to 100 for all experiments.

In this work, all experiments are carried out with PyTorch 1.8 on NVIDIA RTX 3090 GPU in the Ubuntu22.04 environment. Codes will be available after acceptance at <https://github.com/JL-tong/PDTrans.git>.

5.3. Accuracy Analysis

We first evaluate the performance of long-term and short-term forecasting on Electricity and Traffic datasets with the

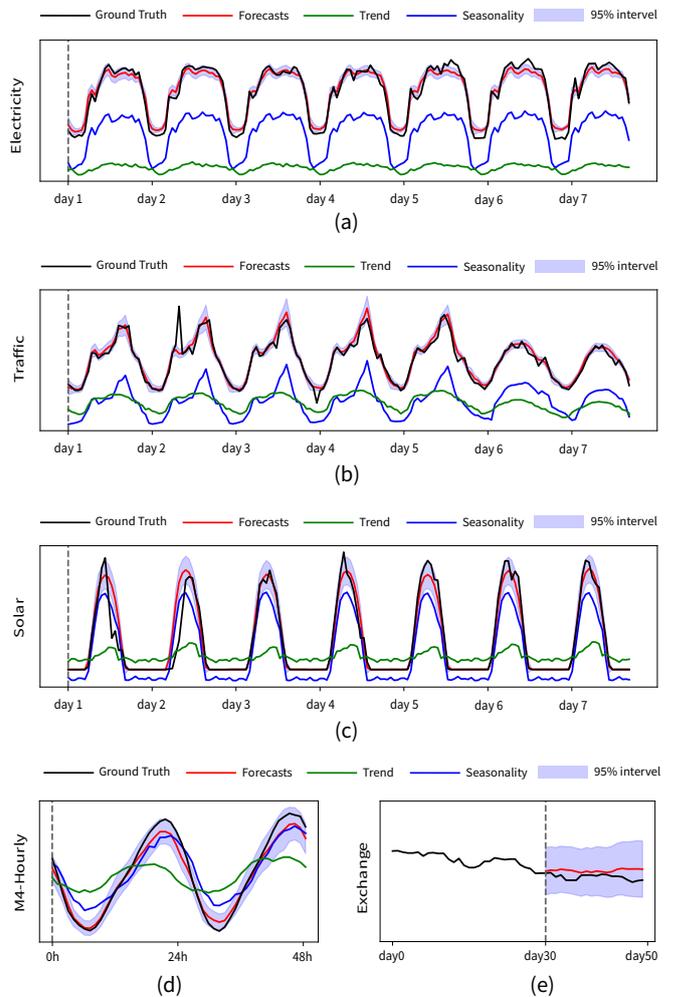


Figure 2: The forecasting results. (a) Result on Electricity. (b) Result on Traffic. (c) Result on Solar. (d) Result on M4-Hourly. (e) Result on Exchange.

Table 5: Ablation study. The performance of $\rho_{0.5}/\rho_{0.9}$ loss. w/ denotes with probabilistic decomposition module and w/o denotes without probabilistic decomposition module.

	Electricity	Traffic	Solar	Exchange	M4-Hourly
PDTrans(w/o)	0.064 / 0.036	0.134 / 0.090	0.214 / 0.086	0.018 / 0.008	0.046 / 0.028
PDTrans(w/)	0.058 / 0.030	0.113 / 0.088	0.205 / 0.073	0.011 / 0.006	0.033 / 0.023

Table 6: Ablation study on Electricity dataset. The effect of the trade-off coefficients β and γ .

	$\rho_{0.5}$	$\rho_{0.9}$
$\gamma = 0.5, \beta = 0.5$	0.058	0.032
$\gamma = 1, \beta = 1$	0.058	0.030
$\gamma = 5, \beta = 5$	0.059	0.031
$\gamma = 10, \beta = 10$	0.064	0.035

Table 7: Comparing the model of DeepAR and the model of DeepAR equipped with probabilistic decomposition.

	Electricity	Traffic
DeepAR	0.075 / 0.040	0.161 / 0.099
PD-DeepAR	0.072 / 0.036	0.147 / 0.096

probabilistic decomposition Transformer model. Both the Electricity and Traffic datasets are divided into 7 days of data for testing. For short-term forecasts, we follow [10] to assess the performance by rolling forecasts for 7 days, with each forecast being 24 hours long. For the long-term forecasting, the length of prediction range is specified as 7 days directly, with 14 days (i.e., 336 observations per time series) for conditioning range. The comparison methods are autoregressive probabilistic models include DeepAR[21], DeepSSM[22], ConvTrans[10], and SSDNet[32]. The result is shown in Table 3.

Our model achieves the best performance for both long-term and short-term forecasting on Electricity dataset. Besides, the PDTrans get better results in long-term forecasting both for $\rho_{0.5}$ and $\rho_{0.9}$ on Traffic dataset. For short-term forecasting on Traffic dataset, the proposed model achieved the best $\rho_{0.5}$ and a competitive $\rho_{0.9}$. The results indicate that the hierarchical forecasting mechanism can effectively keep the long-term dependency of the model and improve forecasting performance.

To more comprehensively assess the effectiveness of the proposed model, the PDTrans is compared with the mainstream methods on five real-world datasets. For the hourly interval dataset, the model input length is set to 7 days and the prediction length is specified as 24 hours. For the day-interval dataset, the model input is set to 30 and the prediction length is selected as 20, following [32]. The comparison methods include Prophet [34], DeepAR[21], DeepSSM[22], ConvTrans[10], N-Beats[35], Informer[31], Autoformer[28], and SSDNet[32]. The $\rho_{0.5}$ and $\rho_{0.9}$ metrics are reported for probabilistic forecasting models. Besides, we only report $\rho_{0.5}$ metrics for the non-probabilistic forecasting models.

For Electricity, Solar, and Exchange datasets, the proposed PDTrans has the best $\rho_{0.5}$ and $\rho_{0.9}$ among the models involved in the comparison. Besides, our method achieves the best perfor-

mance for $\rho_{0.5}$ metrics on Traffic dataset, and a very competitive performance for the $\rho_{0.9}$ metrics. The PDTrans performs better than autoregressive probabilistic forecasting models, such as DeepAR and ConvTrans, which indicates that hierarchical probabilistic forecasting can improve the performance of autoregressive probabilistic forecasting. In addition, our model performs well on datasets with significant periodicities such as the Electricity and Solar datasets, indicating that the model has an excellent ability for mining periodic features.

5.4. Interpretability Analysis

The model in this paper extracts the short-term trend instead of the long-term trend since the long-term trend requires obtaining a long series, which is difficult to operate in practice.

Figure 2 presents the results of a rolling 7-day forecasting on the Electricity, Traffic, and Solar datasets, where the black line represents the ground truth, the red line denotes the forecasting result, the green line is the trend term of forecasting result, and the blue line is the seasonality term of forecasting result. The purple shaded area of the figure represents the 95% confidence interval. For the M4-Hourly dataset, only the 48-hour ahead forecasting results are presented in the figure. Besides, the decomposition results of the Exchange dataset are not shown, due to the lack of regular fluctuations.

The results in Figure 2 are interpretable that the trend curve is moving proposed and reflects the trend of change, the seasonality term presents regular fluctuations. For instance, in Figure 2 (b) from day 5 to day 7, the green line not only presents the daily trend, but also reflects the decreasing trend of the peak.

5.5. Ablation Study

We conducted ablation experiments to evaluate the effectiveness of the probabilistic decomposition module part of PDTrans. Specifically speaking, we compare the PDTrans with the model eliminating the probability decomposition module that is equivalent to the autoregressive probabilistic Transformer. The short-term forecasting performance of the two models on the five datasets is shown in Table 5. For $\rho_{0.5}$ metrics, the PDTrans decreases by 9%, 15%, 4%, 38%, 28% for the Electricity, Traffic, Solar, Exchange, and M4-Hourly, respectively. A similar phenomenon can be observed for the $\rho_{0.9}$ metrics. The results indicate that the model equipped with the probabilistic decomposition module achieves better performance.

In addition, to analyze the effect of different trade-off coefficients on the model, we designed relevant experiments on the Electricity dataset. The result can be found in Table 6. The experimental results show that the model performance fluctuates slightly when the coefficients fluctuate within a certain range. For example, when the coefficients β and γ both change from 1

to 5, the $\rho_{0.5}$ evaluation indicator increases by only 1.69% from 0.058 to 0.059 and the $\rho_{0.9}$ increases by 3.23% from 0.030 to 0.031. The results show that the PDTrans model is not sensitive to the trade-off coefficients within a certain range, demonstrating satisfactory robustness.

5.6. Further Exploration

To further explore the benefits of the probabilistic decomposition module to time series forecasting tasks, the probabilistic decomposition module is employed to the DeepAR [21], which is an autoregressive probabilistic forecasting network based on LSTM. Similar to PDTrans, the probabilistic decomposition module is employed to model the likelihood that output by LSTM for hierarchical probabilistic and interpretable forecasting. The forecasting result on Electricity and Traffic datasets can be found in Table 7. The forecasting result shown in Table 7 indicates that the PD-DeepAR performs significantly better than DeepAR on both datasets. It indicates that the probability decomposition module can effectively improve the forecasting performance of the model through hierarchical forecasting.

6. Conclusion

In this work, we propose probabilistic decomposition Transformer (PDTrans) model for hierarchical and interpretable probabilistic forecasting of intricate time series data, where the PDTrans consists of Transformer and conditional generative model. The Transformer is employed to extract temporal patterns and implement primary autoregressive probabilistic forecasting. The conditional generative model achieves hierarchical probabilistic forecasting through variational inference, which can effectively reduce the impact of exposure bias in the autoregressive forecasting process. In addition, the conditional generative model generates trends and seasonality features by probabilistic decoders to achieve separation of intricate patterns and interpretable forecasts. A series of ablation experiments are designed to demonstrate the effectiveness and robustness of the probabilistic decomposition block. Moreover, the performance of the PDTrans is evaluated on five time series datasets, showing that it compares favorably with state of the art in terms of accuracy. The results indicate that the PDTrans is a reliable alternative to time series forecasting.

Acknowledgements

This work was supported in part by the Zhishan Young Scholar Program of Southeast University; in part by the Fundamental Research Funds for the Central Universities under Grant 2242021R41118. Besides, we thank the Big Data Computing Center of Southeast University for providing the facility support on the numerical calculations.

Appendix A. Evidence Lower Bound of Log-Likelihood

The derivation for the evidence lower bound of log-likelihood is denoted as follows:

$$\begin{aligned}
& \log p_{\theta}(\mu_{t_0+1:t_0+\tau} | Y_{1:t_0}) \\
&= D_{KL}(q_{\phi}(z | Y_{1:t_0}, \mu_{t_0+1:t_0+\tau}) || p_{\theta}(z | Y_{1:t_0}, \mu_{t_0+1:t_0+\tau})) \\
&\quad + \mathbb{E}_{q_{\phi}(z | Y_{1:t_0}, \mu_{t_0+1:t_0+\tau})} [- \log q_{\phi}(z | Y_{1:t_0}, \mu_{t_0+1:t_0+\tau})] \\
&\quad + \mathbb{E}_{q_{\phi}(z | Y_{1:t_0}, \mu_{t_0+1:t_0+\tau})} [\log p_{\theta}(\mu_{t_0+1:t_0+\tau}, z | Y_{1:t_0})] \\
&\geq \mathbb{E}_{q_{\phi}(z | Y_{1:t_0}, \mu_{t_0+1:t_0+\tau})} [- \log q_{\phi}(z | Y_{1:t_0}, \mu_{t_0+1:t_0+\tau})] \\
&\quad + \mathbb{E}_{q_{\phi}(z | Y_{1:t_0}, \mu_{t_0+1:t_0+\tau})} [\log p_{\theta}(\mu_{t_0+1:t_0+\tau}, z | Y_{1:t_0})] \quad (\text{A.1}) \\
&= \mathbb{E}_{q_{\phi}(z | Y_{1:t_0}, \mu_{t_0+1:t_0+\tau})} [- \log q_{\phi}(z | Y_{1:t_0}, \mu_{t_0+1:t_0+\tau})] \\
&\quad + \mathbb{E}_{q_{\phi}(z | Y_{1:t_0}, \mu_{t_0+1:t_0+\tau})} [\log p_{\theta}(z | Y_{1:t_0})] \\
&\quad + \mathbb{E}_{q_{\phi}(z | Y_{1:t_0}, \mu_{t_0+1:t_0+\tau})} [\log p_{\theta}(\mu_{t_0+1:t_0+\tau} | Y_{1:t_0}, z)] \\
&= - D_{KL}(q_{\phi}(z | Y_{1:t_0}, \mu_{t_0+1:t_0+\tau}) || p_{\theta}(z | Y_{1:t_0})) + \\
&\quad \mathbb{E}_{q_{\phi}(z | Y_{1:t_0}, \mu_{t_0+1:t_0+\tau})} [\log p_{\theta}(\mu_{t_0+1:t_0+\tau} | z, Y_{1:t_0})].
\end{aligned}$$

References

- [1] Zi, W., Xiong, W., Chen, H. Chen, L. TAGCN: Station-level demand prediction for bike-sharing system via a temporal attention graph convolution network. *Information Sciences*. 561 (2021), pp. 274-285.
- [2] Tong, J., Xie, L., Fang, S., Yang, W. Zhang, K. Hourly solar irradiance forecasting based on encoder-decoder model using series decomposition and dynamic error compensation. *Energy Conversion And Management*. 270 (2022), pp. 116049.
- [3] Ali, A., Zhu, Y. Zakarya, M. Exploiting dynamic spatio-temporal correlations for citywide traffic flow prediction using attention based neural networks. *Information Sciences*. 577 (2021), pp. 852-870.
- [4] Sun, J., Fujita, H., Zheng, Y. Ai, W. Multi-class financial distress prediction based on support vector machines integrated with the decomposition and fusion methods. *Information Sciences*. 559 (2021), pp. 153-170.
- [5] Box, G., Jenkins, G., Reinsel, G. Ljung, G. Time series analysis: forecasting and control. *John Wiley & Sons*. (1970).
- [6] Hyndman, R., Koehler, A., Snyder, R. Grose, S. A state space framework for automatic forecasting using exponential smoothing methods. *International Journal Of Forecasting*. 18 (2002), pp. 439-454.
- [7] Chung, Junyoung, et al. A recurrent latent variable model for sequential data. *Advances In Neural Information Processing Systems*. 28 (2015).
- [8] Fraccaro, M., Sønderby, S., Paquet, U. Winther, O. Sequential neural models with stochastic layers. *Advances In Neural Information Processing Systems (NeurIPS)*. (2016).
- [9] Vaswani, A., Shazeer, N., Parmar, N., Uszkoreit, J., Jones, L., Gomez, A., Kaiser, Ł. Polosukhin, I. Attention is all you need. *Advances In Neural Information Processing Systems (NeurIPS)*. (2017).
- [10] Li, S., Jin, X., Xuan, Y., Zhou, X., Chen, W., Wang, Y. Yan, X. Enhancing the locality and breaking the memory bottleneck of transformer on time series forecasting. *Advances In Neural Information Processing Systems (NeurIPS)*. (2019).
- [11] Williams, R. Zipser, D. A learning algorithm for continually running fully recurrent neural networks. *Neural Computation*. 1 (1989), pp. 270-280.
- [12] Lamb, A., Alias Parth Goyal, A., Zhang, Y., Zhang, S., Courville, A. Bengio, Y. Professor forcing: A new algorithm for training recurrent networks. *Advances In Neural Information Processing Systems*. 29 (2016).
- [13] Bengio, S., Vinyals, O., Jaitly, N. Shazeer, N. Scheduled sampling for sequence prediction with recurrent neural networks. *Advances In Neural Information Processing Systems (NeurIPS)*. (2015).
- [14] Hyndman, R. Athanasopoulos, G. Forecasting: principles and practice. *OTexts*. (2018).
- [15] Cleveland, R., Cleveland, W., McRae, J. Terpenning, I. STL: A seasonal-trend decomposition. *J. Off. Stat.* 6 (1990), pp. 3-73.

- [16] Kingma, D. Welling, M. Auto-encoding variational bayes. *ArXiv Preprint ArXiv:1312.6114*. (2013).
- [17] Rezende, D., Mohamed, S. Wierstra, D. Stochastic backpropagation and approximate inference in deep generative models. *International Conference On Machine Learning (ICML)*. (2014), pp. 1278-1286.
- [18] Billah, B., King, M., Snyder, R. Koehler, A. Exponential smoothing model selection for forecasting. *International Journal Of Forecasting*. 22 (2006), pp. 239-247.
- [19] Wen, R., Torkkola, K., Narayanaswamy, B. Madeka, D. A multi-horizon quantile recurrent forecaster. *ArXiv Preprint ArXiv:1711.11053*. (2017).
- [20] Hewamalage, H., Bergmeir, C. Bandara, K. Recurrent neural networks for time series forecasting: Current status and future directions. *International Journal Of Forecasting*. 37 (2021), pp. 388-427.
- [21] Salinas, D., Flunkert, V., Gasthaus, J. Januschowski, T. DeepAR: Probabilistic forecasting with autoregressive recurrent networks. *International Journal Of Forecasting*. 36 (2020), pp. 1181-1191.
- [22] Rangapuram, S., Seeger, M., Gasthaus, J., Stella, L., Wang, Y. Januschowski, T. Deep state space models for time series forecasting. *Advances In Neural Information Processing Systems (NeurIPS)*. (2018).
- [23] Lai, G., Chang, W., Yang, Y. Liu, H. Modeling long-and short-term temporal patterns with deep neural networks. *ACM SIGIR Conference On Research And Development In Information Retrieval (SIGIR)*. (2018).
- [24] Bai, S., Kolter, J. Koltun, V. An empirical evaluation of generic convolutional and recurrent networks for sequence modeling. *ArXiv Preprint ArXiv:1803.01271*. (2018).
- [25] Sen, R., Yu, H. Dhillon, I. Think globally, act locally: A deep neural network approach to high-dimensional time series forecasting. *Advances In Neural Information Processing Systems (NeurIPS)*. (2019).
- [26] Chen, Y., Kang, Y., Chen, Y. Wang, Z. Probabilistic forecasting with temporal convolutional neural network. *Neurocomputing*. 399 (2020), pp. 491-501.
- [27] Wen, Q., Zhou, T., Zhang, C., Chen, W., Ma, Z., Yan, J. Sun, L. Transformers in time series: A survey. *ArXiv Preprint ArXiv:2202.07125*. (2022).
- [28] Wu, H., Xu, J., Wang, J. Long, M. Autoformer: Decomposition transformers with auto-correlation for long-term series forecasting. *Advances In Neural Information Processing Systems (NeurIPS)*. (2021).
- [29] Zhou, T., Ma, Z., Wen, Q., Wang, X., Sun, L. Jin, R. FEDformer: Frequency enhanced decomposed transformer for long-term series forecasting. *ArXiv Preprint ArXiv:2201.12740*. (2022).
- [30] Liu, S., Yu, H., Liao, C., Li, J., Lin, W., Liu, A. Dustdar, S. Pyraformer: Low-complexity pyramidal attention for long-range time series modeling and forecasting. *International Conference On Learning Representations (ICLR)*. (2022).
- [31] Zhou, H., Zhang, S., Peng, J., Zhang, S., Li, J., Xiong, H. Zhang, W. Informer: Beyond efficient transformer for long sequence time-series forecasting. *Proceedings Of The AAAI Conference On Artificial Intelligence*. (2021).
- [32] Lin, Y., Koprinska, I. Rana, M. SSDNet: State space decomposition neural network for time series forecasting. *2021 IEEE International Conference On Data Mining (ICDM)*. (2021).
- [33] West, M. Time series decomposition. *Biometrika*. 84 (1997), pp. 489-494.
- [34] Taylor, S. Letham, B. Forecasting at scale. *The American Statistician*. 72 (2018), pp. 37-45.
- [35] Oreshkin, B., Carpov, D., Chapados, N. Bengio, Y. N-BEATS: Neural basis expansion analysis for interpretable time series forecasting. *International Conference On Learning Representations (ICLR)*. (2019).
- [36] Nguyen, N. Quanz, B. Temporal latent auto-encoder: A method for probabilistic multivariate time series forecasting. *Proceedings Of The AAAI Conference On Artificial Intelligence*. (2021).
- [37] Zhang, Y., Zhou, B., Cai, X., Guo, W., Ding, X. Yuan, X. Missing value imputation in multivariate time series with end-to-end generative adversarial networks. *Information Sciences*. 551 (2021), pp. 67-82.
- [38] Miao, X., Wu, Y., Wang, J., Gao, Y., Mao, X. Yin, J. Generative semi-supervised learning for multivariate time series imputation. *Proceedings Of The AAAI Conference On Artificial Intelligence*. 35 (2021), 8983-8991.
- [39] Yoon, J., Jarrett, D. Schaar, M. Time-series generative adversarial networks. *Advances In Neural Information Processing Systems(NeurIPS)*. (2019).
- [40] Jha, P., Bohlke-Schneider, M., Mercado, P., Kapoor, S., Nirwan, R., Flunkert, V., Gasthaus, J. Januschowski, T. PSA-GAN: Progressive Self Attention GANs for Synthetic Time Series. *International Conference On Learning Representations (ICLR)*. (2021).
- [41] Han, W., Wang, L., Feng, R., Gao, L., Chen, X., Deng, Z., Chen, J. Liu, P. Sample generation based on a supervised Wasserstein Generative Adversarial Network for high-resolution remote-sensing scene classification. *Information Sciences*. 539 (2020), pp. 177-194.
- [42] Sohn, K., Lee, H. Yan, X. Learning structured output representation using deep conditional generative models. *Advances In Neural Information Processing Systems (NeurIPS)*. (2015).
- [43] Fortuin, V., Baranchuk, D., Rättsch, G. Mandt, S. Gp-vae: Deep probabilistic time series imputation. *International Conference On Artificial Intelligence And Statistics*. (2020), pp. 1651-1661.
- [44] Tang, B. Matteson, D. Probabilistic transformer for time series analysis. *Advances In Neural Information Processing Systems (NeurIPS)*. (2021).
- [45] Rasul, K., Seward, C., Schuster, I. Vollgraf, R. Autoregressive denoising diffusion models for multivariate probabilistic time series forecasting. *International Conference On Machine Learning (ICML)*. (2021), pp. 8857-8868.
- [46] Li, J., Wu, B., Sun, X. Wang, Y. Causal Hidden Markov Model for Time Series Disease Forecasting. *Proceedings Of The IEEE/CVF Conference On Computer Vision And Pattern Recognition (CVPR)*. (2021), pp. 12105-12114.
- [47] Devlin, J., Chang, M., Lee, K. Toutanova, K. Bert: Pre-training of deep bidirectional transformers for language understanding. *ArXiv Preprint ArXiv:1810.04805*. (2018).
- [48] Hochreiter, S. Schmidhuber, J. Long short-term memory. *Neural Computation*. 9 (1997), pp. 1735-1780.
- [49] Cho, K., Van Merriënboer, B., Gulcehre, C., Bahdanau, D., Bougares, F., Schwenk, H. Bengio, Y. Learning phrase representations using RNN encoder-decoder for statistical machine translation. *Conference On Empirical Methods In Natural Language Processing (EMNLP)*. (2014).
- [50] Higgins, I., Matthey, L., Pal, A., Burgess, C., Glorot, X., Botvinick, M., Mohamed, S. Lerchner, A. beta-vae: Learning basic visual concepts with a constrained variational framework. *International Conference On Learning Representations (ICLR)*. (2016).
- [51] Kingma, D., Salimans, T. Welling, M. Variational dropout and the local reparameterization trick. *Advances In Neural Information Processing Systems (NeurIPS)*. (2015).
- [52] Makridakis, S., Spiliotis, E. Assimakopoulos, V. The M4 Competition: Results, findings, conclusion and way forward. *International Journal Of Forecasting*. 34 (2018), pp. 802-808.
- [53] Kingma, Diederik P and Ba, Jimmy. Adam: A method for stochastic optimization. *In International Conference on Learning Representations (ICLR)*. (2015).