Multivariate Probabilistic Time Series Forecasting via Conditioned Normalizing Flows

Kashif Rasul¹, Abdul-Saboor Sheikh, Ingmar Schuster, Urs Bergmann and Roland Vollgraf

Zalando Research @ ICLR 2021

'**Email:**kashif.rasul@zalando.de



1/12

Introduction

- Time series forecasting ubiquitous in academia and industry
- Classical methods do not scale with large data
- Deep learning methods offer a popular alternative -

Probabilistic Forecasts

- Probabilistic forecasts account for uncertainty
- Individual time series can be statistically dependent
- Model needs to account for this via multivariate forecasts
 - → e.g. in traffic flow a disruption will ripple to nearby streets etc.

This work proposes a general deep learning based probabilistic multivariate time series forecasting model.



	fest		
		7	
と			
		Ł	\rightarrow

Univariate Problem

- * $x_t^i \in \mathbb{R}$ where $i \in \{1, \dots, D\}$ and t time index
- → Context window $[1, t_0)$ and prediction window $[t_0, T]$
- $arrow ext{Covariates for each series at } t ext{ given by vector } \mathbf{c}_t^i \in \mathbb{R}^C$
- Univariate probablistic problem:

$$p_{\mathcal{X}}(x^i_{t_0:T}|x^i_{1:t_0-1},\mathbf{c}^i_{1:T}) pprox \Pi^T_{t=t_0} p_{ heta}(x^i_t|x^i_{1:T})$$

dex $\textit{vindow}\left[t_0,T ight]$ J $\textit{vector}\,\mathbf{c}_t^i \in \mathbb{R}^C$

 $_{:t-1}, \mathbf{c}_{1:T}^i)$



Multivariate Forecast?

- \Rightarrow Can use the same model with vector valued input $\mathbf{x}_t \in \mathbb{R}^D$
- Covariates now associated with the whole vector rather than individual entities
- → Need a full joint distribution model for each t
- → E.g. multivariate Gaussian but:
 - ⇒ increase parameters by $O(D^2)$
 - ⇒ log-prob will be $O(D^3)$

Put $\mathbf{x}_t \in \mathbb{R}^D$: tor rather than

Normalizing Flows (Rezende & Mohamed, 2015)

- → NF are a sequence of invertible mappings from $\mathbb{R}^D \mapsto \mathbb{R}^D$
- → Idea: transform data distribution to some simple distribution and max. that log-prob or sample from simple and go backwards
- → Has two properties:
 - I. inverse is easy to evaluate
 - 2. Jacobian determinant is O(D)



Real NVP (Dinh et al., 2017)

- → Introduces a bijection: Coupling Layer
- → For some *d* and two neural networks $s_{\theta}: \mathbb{R}^{d} \to \mathbb{R}^{D-d}$ and $t_{\theta}: \mathbb{R}^{d} \to \mathbb{R}^{D-d}$ we have: $\begin{cases} \mathbf{y}^{1:d} = \mathbf{x}^{1:d} \\ \mathbf{y}^{d+1:D} = \mathbf{x}^{d+1:D} \odot \exp(s_{\theta}(\mathbf{x}^{1:d})) + t_{\theta}(\mathbf{x}^{1:d}) \end{cases}$
- Jacobian is block triangular: det. is sum over diagonal
- → Conditioning: concat vectors to inputs of s_{θ} and t_{θ}





11/12

Summary

- Experiments on 6 real-world datasets: establish SOTA results (with max. D = 2,000)
- Extensive comparison against competitive modern probablistic multivariate methods
- Can replace RNN with self-attention based Transformer (Vaswani et al., 2017) module
- PyTorch implementation Github: zalandoresearch/pytorch-ts

12/12