Decomposition-based Data Augmentation for Time-series Building Load Data

Yang Deng, Rui Liang, Dan Wang, Ao Li and Fu Xiao

The Hong Kong Polytechnic University





Department of Building Environment and Energy Engineering 建築環境及能源工程學系

Background



Time-series building load data

- The fundamental of many applications, such like load forecasting, etc.
- Load data is collected gradually over time.



Building data Augmentation

- □ ML-based applications require a large amounts of data for training/testing ML model.
- Data augmentation (DA) scheme is necessary.

Existing Data Augmentation scheme



Most are based on generative models, i.e., GANs.



Existing Data Augmentation scheme A long period, usually year-level Step1: Collection **Target building** Januarv JUCH J. generation Step 4: 1.1

Support

applications

4

Existing Data Augmentation scheme



A long period, usually year-level





The problem under insufficient time coverage



The problem under insufficient time coverage













Potential Approach:

Data Augmentation based on *Time-series decomposition*

Time-series decomposition



Time series decomposition assumes that the time-series can be regarded as a collection of several components.



[1] M. West. 1997. Time series decomposition. Biometrika (1997).

Decomposition-based Data augmentation 🍪

The rationale:

Augmenting the data of each component can be a more targeted process.



Trend ? Seasonality ? Irregular ? Analyzed data:
407 USA buildings from Genome dataset [2]

[2] C. Miller, A. Kathirgamanathan, et al. 2020. The building data genome project 2, energy meter data from the ASHRAE great energy predictor III competition. Scientific data (2020).

Trend Seasonality Irregular - e.g., the rate at which equipment is aging.

13



Does these components exist in building load data? Trend Seasonality

Irregular Reflects noises for a short duration, i.e., spike/dip loads within a stable load period







Data Augmentation scheme Design

Design Overview



Goal

- To minimize the distance between the synthetic load and the real data, given the target building with an insufficient distribution of data collection.
- Decomposition-based augmentation scheme (DAST)





Step 1: Load time-series decomposition

- Challenge 1: Decompose the sequence of { D } into { C(p) } and { R }
- Solution 1: A classical season-trend decomposition (STL) method, with a calibration for day type.







Challenge 2: Minimize the distance between the augmented component and the real ones.

Step 2: Augmentation for the components

(p)

Solution:

- A k-means clustering to recognize the daily patterns.
- A classical Pattern mixing method to generate daily patterns.

Solution:

- A learning-based domain translation algorithm to learn transformation operations.
- A metadata-clustering method to search the training data.

Solution:

 A statistics-based method (KDE)

Challenge 2: Minimize the distance of the augmented component and the real ones.



Challenge 3: To remove the low-fidelity D'

Step 3: combination

Ì

Solution:

• A contrastive learning-based module to learn the relation of *C(p)* and *D*, and then recognize the low-fidelity *D'* based on binary-classification.



 $D'(\hat{C}(\hat{p}), \hat{R})$ Synthetic Daily load

Challenge 3: To remove the low-fidelity D'

Evaluation Setup



Dataset

- Six buildings according to different dimensions:
 - building type, location, and building size.
- Settings of given data size
 - Two weeks / one month / three months
- Baselines:
 - RCGAN
 - RCGAN adopts RNN in Vanilla GAN, a benchmark model for time-series data augmentation.
 - □ <u>cVAE</u>
 - conditional variational autoencoder. It can generate data for specific types of months.
 - TimeGAN
 - A state-of-the-art GAN for time-series scenarios.

Building	The ID in Genome	Location	Floor area (sqft)
A	Rat_public_Emilee	Washington	22500
В	Rat_public_Isabel	Washington	16576
С	Fox_office_Joy	Tempe	70837
D	Fox_education_Virginia	Tempe	12773
E	Bear_education_Iris	Berkeley	58733
F	Peacock_education_Ophelia	Princeton	120836

Improvement of DAST

Quantitative evaluation:

DAST

3.0

2.5

2.0

1.5

1.0

0.5

0.0

MMD

 MMD (maximum mean discrepancy), lower the better

ZZ CVAE

B

data and synthetic data) by 49%.

Reduce the distance (distribution between real

venient of DAST

TimeGAN

RCGAN





Qualitative evaluation:

Case study: the benefit to building load forecasting (BLF) tasks

- Task 1: BLF model training
 - Use DAST to support RNN, LSTM training
 - Train-on-synthetic and test-on-real (TSTR) I



- Task 2: Trained BLF model testing
 - Differentiate the performance of 30 BLF models for a target building.
 - □ Top-5 ranking.



Conclusion



- 1. We proposed DAST, a decomposition-based data augmentation scheme for insufficient data distribution scenario.
- 2. We analyzed decomposed components in real buildings and developed appropriate augmentation schemes.
- 3. We conducted qualitative and quantitative experiments to assess the quality of generated data.
- 4. We conducted case study on building load forecasting tasks.



Thanks for listening!