AIM 2: Artificial Intelligence in Medicine II Harvard - BMIF 203 and BMI 702, Spring 2025

Lecture 13: Digital biomarkers and disease progression tracking, Patient/disease progression modeling using transformers, Time series in healthcare, Longitudinal EHR modeling



For the Study of Natural & Artificial Intelligence at Harvard University



Marinka Zitnik marinka@hms.harvard.edu

Time series are everywhere

Irregularly sampled with varying time intervals between successive readouts, complex dynamics and various sensors observed at different time points

Climate





Healthcare





Space systems





Domain Adaptation for Time Series Under Feature and Label Shifts, ICML 2023 Graph-Guided Network for Irregularly Sampled Multivariate Time Series, ICLR 2022

Self-Supervised Contrastive Pre-Training For Time Series via Time-Frequency Consistency, NeurIPS 2022 Encoding Time-Series Explanations through Self-Supervised Model Behavior Consistency, NeurIPS 2023

Time series in healthcare



Mihaela van der Schaar

Time series in healthcare

- Multiple streams of measurements
- Measurements are sparse, irregularly and informatively sampled
- Multiple outcomes of interest (various events of interest, various morbidities)
- True clinical states are unobserved (e.g., onset of diseases)
- Many possible patterns (heterogeneous phenotypes, comorbidities)



Time series data come in different forms

- Time series consist of
 - Data points ordered in time
 - Preferably in regular intervals
 - Typically, 1 or more components/channels/dimensions



- Time axis not strictly needed but helpful to
 - Order data
 - Perform analysis and forecast future readouts
 - Temporal resolution of the time series depends on the use-case
 - Often in minutes, hours, days, weeks, month, ...

Today's lecture

Time series tasks in healthcare

2 Self-supervised pre-training for time series

3 Learning representations of regular and irregular time series

Understanding time series models

Time series tasks in healthcare **Dynamic forecasting**

Time-to-event and survival analysis

Clustering and phenotyping

Screening and monitoring

Early prediction of diagnosis

Treatment effect prediction

Time series forecasting: Predicting the future

Univariate forecasting

- Single time series forecasting (with time)
- The future may be forecasted just by looking at the past
- Analysis and ML forecasting methods for this are widely researched
- Simple methods like moving average, regression, and general neural networks can often be sufficient, depending on the accuracy required



Multivariate forecasting

- Forecasting with multiple time series
- Access to conditional past and future data, both continuous and categorical
- Accurate forecasting is a difficult task
 - Dependency on multiple covariates
 - Long-range dependencies
 - Inherent uncertainty in input/target data



Time series forecasting: Predicting the future

- Build disease progression models
- Understand and model carefully the available data!
- Learn the model parameters from available EHR data (Training time)
- Issue dynamic forecasts for the patient at hand (Test time/Run-time)
- Unravel new understanding of disease progression
 - Population
 - Sub-groups of patients
 - Personalized



Disease progression models: Formalism

Markov models

 $P(\boldsymbol{Z}_{n+1} | \mathcal{H}_{t_n}) = P(\boldsymbol{Z}_{n+1} | \boldsymbol{Z}_n)$



Disadvantages:

- Observable models
- One disease at a time
- "Average" patient

Population-level representation of disease states



Disease progression models: Formalism

Hidden Markov Models (HMMs)

Introducing latent (hidden/unobservable) disease states

Hidden states



Observations

Disease Stages

Clinical findings Lab measurements Vital signs Treatments Events of interest Observation times

Disease progression models: Goals

History matters!

Ignore history

- Previous states
- Order of states
- Duration in a state

One size fits all!

Only capture population-level transitions across progression stages Ignores individual clinical trajectories

Goal A: Accurately forecast individual-level disease trajectories

What are the risks of mortality, relapse, comorbidities, complications, etc. in the future?

Goal B: Understand disease progression mechanisms

Underlying latent structure of disease evolution

Patients' <u>subgroup</u> analysis

Refined phenotypes

Mihaela van der Schaar

Clustering & phenotyping: How should we group patients?

Example of 3 patients diagnosed with breast cancer (BC)

Should we group patients based on similarity in the time-series observations?



Clustering & phenotyping: How should we group patients?

Example of 3 patients diagnosed with breast cancer (BC)

What if both Patient A and C will have an adverse event (e.g., death) that can be expected by increases in cancer antigen and mammographic density



- Predictive of similar future outcomes
- Doctors and patients can actively plan
- Learn representations of past observations (time-series data) that best describe future events and outcomes

Personalized screening and monitoring



Who to screen? When to screen? What to screen?

- What is the value of various information over time for this event for this individual?

Personalized screening: Formalism



- Deep Sensing [Yoon, Jordon, vdS, 2018]
- Disease Atlas [Yoon, Jordon, vdS, 2019]
- Clairvoyance [Jarrett et al, 2021]

Early prediction of diagnosis



Early prediction of diagnosis

Risk prediction

- Segments individuals using population-based risks, usually based on few variables
- Rarely uses longitudinal data usually only calculated once
- Risk scores then lead to guideline-driven management of patients often rigid
 - Many diseases lack guidelines and protocols
 - This is all predicated upon a quantitative understanding of disease progression

How can we detect disease early?

- Early diagnosis is more than just event prediction/forecasting
- It involves unravelling and dissecting the underlying states of disease progression towards the event of interest



A quantitative understanding of disease progression is needed!

Mihaela van der Schaar

Early diagnosis: How?

Healthy Emerging Pre-malignant Conditions (Dysplasia)



Early diagnosis: How?

Healthy

Emerging Pre-malignant Conditions (Dysplasia)



Disease Progresses



Treatment effect prediction



Treatment effect prediction



Causal effect inference from longitudinal patient observational data



Challenges in using longitudinal observational data for estimating treatment outcomes

The patient history $\bar{\mathbf{H}}_t = (\bar{\mathbf{X}}_t, \bar{\mathbf{A}}_{t-1}, \mathbf{V})$ contains time-dependent confounders which bias the treatment assignment \mathbf{A}_t in the observational dataset.

Patient covariates - affected by past treatments which then influence future treatments and outcomes



Bias from time-dependent confounders.

Time series tasks in healthcare **Dynamic forecasting**

Time-to-event and survival analysis

Clustering and phenotyping

Screening and monitoring

Early prediction of diagnosis

Treatment effect prediction

Today's lecture

Time series tasks in healthcare

2 Self-supervised pre-training for time series

3 Learning representations of regular and irregular time series

4 Understanding time series models

Pre-training on time series

 Question: How to process a time series dataset so as to greatly improve generalization to new time series coming from a different dataset



Goal:

- Transfer a model to new datasets without explicit retraining or minimal adaptation
- Resulting performance \geq SoTA model on target dataset alone

Why is this challenging?

Expected performance gains are often not realized:



Do time series datasets represent unique challenges?

- Target datasets are not available in pre-training
 - Pre-trained model must capture a latent shared property that can apply to an unseen dataset
- Need to identify a shared property universal to different time series datasets to enable transfer from pre-training to target datasets



Shapes, edges, texture



Word order, grammar







Time-frequency consistency (TF-C)

Representational Time-Frequency Consistency (TF-C). Let be given a time series sample x_i . Then in a model \mathcal{F} satisfying TF-C, time-based representation z_i^T and frequency-based representation z_i^F learned from x_i by \mathcal{F} , and representations learned from local augmentations of x_i are close together in the latent time-frequency space.



Overview of TF-C approach



Datasets

- SleepEEG: 371,055 univariate brainwaves (100 Hz) collected from 197 individuals. Samples are labeled by 5 sleeping stages
- Epilepsy: Brain activity of 500 subjects recorded by single-channel EEG (174 Hz). Samples are labeled by epilepsy seizures
- FD-A: Vibration from rolling bearing of a mechanical system aiming at fault detection. Every sample has 5,120 timestamps and an indicator for one out of three mechanical device states
- FD-B: Same as FD-A but the rolling bearings are performed in different working conditions (*e.g.*, varying rotational speed)
- HAR: 10,299 9-dimensional samples recording 6 daily activities
- **GESTURE:** 440 accelerometer samples on 8 hand gestures
- ECG: 8,528 single-sensor ECG recordings sorted into four classes based on human physiology
- EMG: 163 EMG samples with 3-class labels implying muscular diseases









Baselines and experimental setup

- Six pre-training methods:
 - TS-SD, TS2vec, CLOCS, Mixing-up, TS-TCC, SimCLR
- To examine utility of pre-training:
 - Non-deep learning KNN model applied directly to fine-tuning datasets
 - Random initialization approach to randomly initialize fine-tuning model
- Performance metrics: accuracy, precision (macro-averaged), recall, F1 score, AUROC, and AUPRC

Results: One-to-one transfer learning

- Pre-train a model on a <u>pre-training dataset</u> and fine-tune the model on <u>one target dataset only</u>
 - HAR → GESTURE: 6 types of human daily activities measured by an 8channel time series → 8 hand gestures measured by 1 channel

7.2% margin over the best baseline

Models	Accuracy	Precision	Recall	F1 score	AUROC	AUPRC
Non-DL (KNN) Random Init.	$\begin{array}{c} 0.6766 {\pm} 0.0000 \\ 0.4219 {\pm} 0.0865 \end{array}$	$\begin{array}{c} 0.6500 {\pm} 0.0000 \\ 0.4751 {\pm} 0.0925 \end{array}$	$\begin{array}{c} 0.6821 {\pm} 0.0000 \\ 0.4963 {\pm} 0.1026 \end{array}$	$\begin{array}{c} 0.6442 {\pm} 0.0000 \\ 0.4886 {\pm} 0.0967 \end{array}$	$\begin{array}{c} 0.8190 {\pm} 0.0000 \\ 0.7129 {\pm} 0.1206 \end{array}$	$\begin{array}{c} 0.5231 {\pm} 0.0000 \\ 0.3358 {\pm} 0.1194 \end{array}$
TS-SD TS2vec CLOCS Mixing-up	$\begin{array}{c} 0.6937 \pm 0.0533 \\ 0.6453 \pm 0.0260 \\ 0.4731 \pm 0.0229 \\ 0.7183 \pm 0.0123 \\ 0.7592 \pm 0.0545 \end{array}$	$\begin{array}{c} 0.6806 \pm 0.0496 \\ 0.6287 \pm 0.0339 \\ 0.4639 \pm 0.0432 \\ 0.7001 \pm 0.0166 \\ 0.7668 \pm 0.0457 \end{array}$	$\begin{array}{c} 0.6883 {\pm} 0.0525 \\ 0.6451 {\pm} 0.0218 \\ 0.4766 {\pm} 0.0266 \\ 0.7183 {\pm} 0.0123 \\ 0.7566 {\pm} 0.0225 \end{array}$	$\begin{array}{c} 0.6785 {\pm} 0.0495 \\ 0.6261 {\pm} 0.0294 \\ 0.4392 {\pm} 0.0198 \\ 0.6991 {\pm} 0.0145 \\ 0.7457 {\pm} 0.9991 \\ \end{array}$	$\begin{array}{c} 0.8708 \pm 0.0305 \\ 0.8890 \pm 0.0054 \\ 0.8161 \pm 0.0068 \\ \textbf{0.9127} \pm \textbf{0.0018} \\ 0.8866 \pm 0.9127 \end{array}$	$\begin{array}{c} 0.6261 \pm 0.0790 \\ 0.6670 \pm 0.0118 \\ 0.4916 \pm 0.0103 \\ 0.7654 \pm 0.0071 \\ 0.7217 \pm 0.0071 \end{array}$
TS-TCC SimCLR TF-C (Ours)	$\begin{array}{c} 0.7593 {\pm} 0.0242 \\ 0.4383 {\pm} 0.0652 \\ \textbf{0.7824} {\pm} \textbf{0.0237} \end{array}$	$\begin{array}{c} 0.7668 \pm 0.0257 \\ 0.4255 \pm 0.1072 \\ \textbf{0.7982} \pm \textbf{0.0496} \end{array}$	$\begin{array}{c} 0.7566 {\pm} 0.0231 \\ 0.4383 {\pm} 0.0652 \\ \textbf{0.8011} {\pm} \textbf{0.0322} \end{array}$	$\begin{array}{c} 0.7457 \pm 0.0210 \\ 0.3713 \pm 0.0919 \\ \textbf{0.7991} \pm \textbf{0.0296} \end{array}$	$\begin{array}{c} 0.8866 {\pm} 0.0040 \\ 0.7721 {\pm} 0.0559 \\ 0.9052 {\pm} 0.0136 \end{array}$	$\begin{array}{c} 0.7217 {\pm} 0.0121 \\ 0.4116 {\pm} 0.0971 \\ \textbf{0.7861} {\pm} \textbf{0.0149} \end{array}$

On average, our TF-C model claims a large margin of 15.4% over all baselines. Further, the strongest baseline varies across scenarios (i.e., TS-TCC in FD-A \rightarrow FD-B; Mixing-up in SleepEEG \rightarrow Epilepsy)

Results: One-to-many transfer learning

- Pre-train a model on a dataset and use it across a broad <u>range of</u> <u>tasks across many target datasets from diverse domains</u>
 - SleepEEG \rightarrow {EPILEPSY, FD-B, GESTURE, EMG}

Scenarios	Models	Accuracy	Precision	Recall	F1 score	AUROC	AUPRC
SleepEeg ↓ Epilepsy	Non-DL (KNN) Random Init.	$\begin{array}{c} 0.8525 {\pm} 0.0000 \\ 0.8983 {\pm} 0.0656 \end{array}$	$\begin{array}{c} 0.8639 {\pm} 0.0000 \\ 0.9213 {\pm} 0.1369 \end{array}$	$\begin{array}{c} 0.6431 {\pm} 0.0000 \\ 0.7447 {\pm} 0.1135 \end{array}$	$\begin{array}{c} 0.6791 {\pm} 0.0000 \\ 0.7959 {\pm} 0.1208 \end{array}$	$\begin{array}{c} 0.6434 {\pm} 0.0000 \\ 0.8578 {\pm} 0.2153 \end{array}$	$\begin{array}{c} 0.6279 {\pm} 0.0000 \\ 0.6489 {\pm} 0.1926 \end{array}$
	TS-SD TS2vec CLOCS Mixing-up TS-TCC SimCLR TF-C (Ours)	$\begin{array}{c} 0.8952 {\pm} 0.0522 \\ 0.9395 {\pm} 0.0044 \\ \textbf{0.9507} {\pm} 0.0027 \\ 0.8021 {\pm} 0.0000 \\ 0.9253 {\pm} 0.0098 \\ 0.9071 {\pm} 0.0344 \\ 0.9495 {\pm} 0.0249 \end{array}$	$\begin{array}{c} 0.8018 \pm 0.2244 \\ 0.9059 \pm 0.0116 \\ 0.9301 \pm 0.0067 \\ 0.4011 \pm 0.0000 \\ 0.9451 \pm 0.0049 \\ 0.9221 \pm 0.0166 \\ \textbf{0.9456} \pm \textbf{0.0108} \end{array}$	$\begin{array}{c} 0.7647 {\pm} 0.1485 \\ 0.9039 {\pm} 0.0118 \\ \textbf{0.9127} {\pm} 0.0165 \\ 0.5000 {\pm} 0.0000 \\ 0.8181 {\pm} 0.0257 \\ 0.7864 {\pm} 0.1071 \\ 0.8908 {\pm} 0.0216 \end{array}$	$\begin{array}{c} 0.7767{\scriptstyle\pm 0.1855}\\ 0.9045{\scriptstyle\pm 0.0067}\\ \textbf{0.9206}{\scriptstyle\pm 0.0066}\\ 0.4451{\scriptstyle\pm 0.0000}\\ 0.8633{\scriptstyle\pm 0.0215}\\ 0.8178{\scriptstyle\pm 0.0998}\\ 0.9149{\scriptstyle\pm 0.0534} \end{array}$	$\begin{array}{c} 0.7677 {\pm} 0.2452 \\ 0.9587 {\pm} 0.0086 \\ 0.9803 {\pm} 0.0023 \\ 0.9743 {\pm} 0.0081 \\ 0.9842 {\pm} 0.0034 \\ 0.9045 {\pm} 0.0539 \\ \textbf{0.9811} {\pm} \textbf{0.0237} \end{array}$	$\begin{array}{c} 0.7940 {\pm} 0.1825 \\ 0.9430 {\pm} 0.0103 \\ 0.9609 {\pm} 0.0116 \\ 0.9618 {\pm} 0.0104 \\ 0.9744 {\pm} 0.0043 \\ 0.9128 {\pm} 0.0205 \\ \textbf{0.9703 {\pm} 0.0199} \end{array}$
SLEEPEEG ↓ FD-B	Non-DL (KNN) Random Init.	$\substack{0.4473 \pm 0.0000 \\ 0.4736 \pm 0.0623}$	$\begin{array}{c} 0.2847 {\pm} 0.0000 \\ 0.4829 {\pm} 0.0529 \end{array}$	$\substack{0.3275 \pm 0.0000 \\ 0.5235 \pm 0.1023}$	$\begin{array}{c} 0.2284 {\pm} 0.0000 \\ 0.4911 {\pm} 0.0590 \end{array}$	$\begin{array}{c} 0.4946 {\pm} 0.0000 \\ 0.7864 {\pm} 0.0349 \end{array}$	$\begin{array}{c} 0.3308 {\pm} 0.0000 \\ 0.7528 {\pm} 0.0254 \end{array}$
	TS-SD TS2vec CLOCS Mixing-up TS-TCC SimCLR TF-C (Ours)	$\begin{array}{c} 0.5566 {\pm} 0.0210 \\ 0.4790 {\pm} 0.0113 \\ 0.4927 {\pm} 0.0310 \\ 0.6789 {\pm} 0.0246 \\ 0.5499 {\pm} 0.0220 \\ 0.4917 {\pm} 0.0437 \\ \textbf{0.6938 {\pm} 0.0231} \end{array}$	$\begin{array}{c} 0.5710 {\pm} 0.0535 \\ 0.4339 {\pm} 0.0092 \\ 0.4824 {\pm} 0.0316 \\ 0.7146 {\pm} 0.0343 \\ 0.5279 {\pm} 0.0293 \\ 0.5446 {\pm} 0.1024 \\ 0.7559 {\pm} 0.0349 \end{array}$	$\begin{array}{c} 0.6054 {\pm} 0.0272 \\ 0.4842 {\pm} 0.0197 \\ 0.5873 {\pm} 0.0387 \\ \textbf{0.7613 {\pm} 0.0198} \\ 0.6396 {\pm} 0.0178 \\ 0.4760 {\pm} 0.0885 \\ \textbf{0.7202 {\pm} 0.0257} \end{array}$	$\begin{array}{c} 0.5703 {\pm} 0.0328 \\ 0.4389 {\pm} 0.0107 \\ 0.4746 {\pm} 0.0485 \\ 0.7273 {\pm} 0.0228 \\ 0.5418 {\pm} 0.0338 \\ 0.4224 {\pm} 0.1138 \\ \textbf{0.7487} {\pm} \textbf{0.0268} \end{array}$	$\begin{array}{c} 0.7196 {\pm} 0.0113 \\ 0.6463 {\pm} 0.0130 \\ 0.6992 {\pm} 0.0099 \\ 0.8209 {\pm} 0.0035 \\ 0.7329 {\pm} 0.0203 \\ 0.6619 {\pm} 0.0219 \\ \textbf{0.8965 {\pm} 0.0135} \end{array}$	$\begin{array}{c} 0.5693 {\scriptstyle \pm 0.0532} \\ 0.4442 {\scriptstyle \pm 0.0162} \\ 0.5501 {\scriptstyle \pm 0.0365} \\ 0.7707 {\scriptstyle \pm 0.0042} \\ 0.5824 {\scriptstyle \pm 0.0468} \\ 0.5009 {\scriptstyle \pm 0.0477} \\ \textbf{0.7871 {\scriptstyle \pm 0.0267}} \end{array}$
SLEEPEEG ↓ GESTURE	Non-DL (KNN) Random Init.	$\substack{0.6833 \pm 0.0000 \\ 0.4219 \pm 0.0629}$	$\begin{array}{c} 0.6501 {\pm} 0.0000 \\ 0.4751 {\pm} 0.0175 \end{array}$	$\substack{0.6833 \pm 0.0000 \\ 0.4963 \pm 0.0679}$	$\begin{array}{c} 0.6443 {\pm} 0.0000 \\ 0.4886 {\pm} 0.0459 \end{array}$	$\begin{array}{c} 0.8190 {\pm} 0.0000 \\ 0.7129 {\pm} 0.0166 \end{array}$	$\begin{array}{c} 0.5232 {\pm} 0.0000 \\ 0.3358 {\pm} 0.1439 \end{array}$
	TS-SD TS2vec CLOCS Mixing-up TS-TCC SimCLR TF-C (Ours)	$\begin{array}{c} 0.6922 {\pm} 0.0444 \\ 0.6917 {\pm} 0.0333 \\ 0.4433 {\pm} 0.0518 \\ 0.6933 {\pm} 0.0231 \\ 0.7188 {\pm} 0.0349 \\ 0.4804 {\pm} 0.0594 \\ \textbf{0.7642 {\pm} 0.0196} \end{array}$	$\begin{array}{c} 0.6698 {\pm} 0.0472 \\ 0.6545 {\pm} 0.0358 \\ 0.4237 {\pm} 0.0794 \\ 0.6719 {\pm} 0.0232 \\ 0.7135 {\pm} 0.0352 \\ 0.5946 {\pm} 0.1623 \\ \textbf{0.7731 {\pm} 0.0355} \end{array}$	$\begin{array}{c} 0.6867 {\pm} 0.0488 \\ 0.6854 {\pm} 0.0349 \\ 0.4433 {\pm} 0.0518 \\ 0.6933 {\pm} 0.0231 \\ 0.7167 {\pm} 0.0373 \\ 0.5411 {\pm} 0.1946 \\ \textbf{0.7429} {\pm} 0.0268 \end{array}$	$\begin{array}{c} 0.6656 {\pm} 0.0443 \\ 0.6570 {\pm} 0.0392 \\ 0.4014 {\pm} 0.0602 \\ 0.6497 {\pm} 0.0306 \\ 0.6984 {\pm} 0.0360 \\ 0.4955 {\pm} 0.1870 \\ \textbf{0.7572 {\pm} 0.0311} \end{array}$	$\begin{array}{c} 0.8725 {\pm} 0.0324 \\ 0.8968 {\pm} 0.0123 \\ 0.8073 {\pm} 0.0109 \\ 0.8915 {\pm} 0.0261 \\ 0.9099 {\pm} 0.0085 \\ 0.8131 {\pm} 0.0521 \\ \textbf{0.9238 {\pm} 0.0159} \end{array}$	$\begin{array}{c} 0.6185 {\pm} 0.0966 \\ 0.6989 {\pm} 0.0346 \\ 0.4460 {\pm} 0.0384 \\ 0.7279 {\pm} 0.0558 \\ 0.7675 {\pm} 0.0201 \\ 0.5076 {\pm} 0.1588 \\ \textbf{0.7961 {\pm} 0.0109} \end{array}$
	Non-DL (KNN) Random Init.	$\begin{array}{c} 0.4390 {\pm} 0.0000 \\ 0.7780 {\pm} 0.0729 \end{array}$	$\begin{array}{c} 0.3772 {\scriptstyle \pm 0.0000} \\ 0.5909 {\scriptstyle \pm 0.0625} \end{array}$	$\substack{0.5143 \pm 0.0000 \\ 0.6667 \pm 0.0135}$	$\substack{0.3979 \pm 0.0000 \\ 0.6238 \pm 0.0267}$	$\begin{array}{c} 0.6025 {\scriptstyle \pm 0.0000} \\ 0.9109 {\scriptstyle \pm 0.1239} \end{array}$	$\substack{0.4084 \pm 0.0000 \\ 0.7771 \pm 0.1427}$
SleepEeg ↓ Emg	TS-SD TS2vec CLOCS Mixing-up TS-TCC SimCLR TF-C (Ours)	$\begin{array}{c} 0.4606 {\pm} 0.0000 \\ 0.7854 {\pm} 0.0318 \\ 0.6985 {\pm} 0.0323 \\ 0.3024 {\pm} 0.0534 \\ 0.7889 {\pm} 0.0192 \\ 0.6146 {\pm} 0.0582 \\ 0.8171 {\pm} 0.0287 \end{array}$	$\begin{array}{c} 0.1545 {\pm} 0.0000 \\ \textbf{0.8040} {\pm} 0.0750 \\ 0.5306 {\pm} 0.0750 \\ 0.1099 {\pm} 0.0126 \\ 0.5851 {\pm} 0.0974 \\ 0.5361 {\pm} 0.1724 \\ 0.7265 {\pm} 0.0333 \end{array}$	$\begin{array}{c} 0.3333 \pm 0.0000\\ 0.6785 \pm 0.0396\\ 0.5354 \pm 0.0291\\ 0.2583 \pm 0.0456\\ 0.6310 \pm 0.0991\\ 0.4990 \pm 0.1214\\ \textbf{0.8159} \pm 0.0289\end{array}$	0.2111±0.0000 0.6766±0.0501 0.5139±0.0409 0.1541±0.0204 0.5904±0.0952 0.4708±0.1486 0.7683±0.031	$\begin{array}{c} 0.5005 \pm 0.0126\\ \textbf{0.9331} \pm \textbf{0.0164}\\ 0.7923 \pm 0.0573\\ 0.4506 \pm 0.1718\\ 0.8851 \pm 0.0113\\ 0.7799 \pm 0.1344\\ 0.9152 \pm 0.0211 \end{array}$	$\begin{array}{c} 0.3775 {\pm} 0.0110 \\ \textbf{0.8436} {\pm} 0.0372 \\ 0.6484 {\pm} 0.0680 \\ 0.3660 {\pm} 0.1635 \\ 0.7939 {\pm} 0.0386 \\ 0.6392 {\pm} 0.1596 \\ 0.8329 {\pm} 0.0137 \end{array}$

- Pre-training and fine-tuning datasets
 SleepEEG → {FD-B, GESTURE, EMG} are diverse. This gap leads to poor baseline performance
- TF-C has high tolerance to diverse datasets and can can serve as a universal model when no relevant pre-training datasets are available
- TF-C is the best performer in 14/18 settings with an 8.4% performance gain
- TF-C consistently outperforms KNN and Random Init. by a large margin of 42.8% and 25.1%
Is TF-C principle needed for pre-training time series models?

- Removing L_{TF-C}, L_T, and L_F result in performance degradation of 6.1%, 7.2%, and 6.7%
- Replacing \mathcal{L}_{TF-C} with a loss measuring consistency in time (\mathcal{L}_{TT-C}) or frequency (\mathcal{L}_{FF-C}): performance gain cannot be simply explained by contrastive loss; time-frequency consistency is crucial
 - \mathcal{L}_{TF-C} outperforms \mathcal{L}_{TT-C} and \mathcal{L}_{FF-C} by 5.3% and 7.2%

	Accuracy	Precision	Recall	F1 score
W/o \mathcal{L}_{C} and \mathcal{L}_{T}	0.7159+-0.0128	0.7211+-0.0428	0.7246+-0.0428	0.7239+-0.0429
W/o \mathcal{L}_{C} and \mathcal{L}_{F}	0.7327+-0.0328	0.7246+-0.0311	0.7339+-0.0307	0.7317+-0.0356
W/o $\mathcal{L}_{ ext{C}}$	0.7428+-0.0297	0.7289+-0.0278	0.7451+-0.0263	0.7377+-0.0308
Replace \mathcal{L}_{C} with \mathcal{L}_{FF-C}	0.7259+-0.0072	0.7319+-0.0256	0.7338+-0.0133	0.7341+-0.0194
Replace \mathcal{L}_{C} with \mathcal{L}_{TT-C}	0.7124+-0.0091	0.7256+-0.0169	0.7231+-0.0197	0.7296+-0.0209
Full Model (TF-C)	0.7642+-0.0196	0.7731+-0.0355	0.7429+-0.0268	0.7572+-0.0311

Self-Supervised Contrastive Pre-Training For Time Series via Time-Frequency Consistency, NeurIPS 2022

Self-supervised pre-training: Recap

- Pre-training for time series poses unique challenges due to potential mismatches between pre-training and target domains, e.g., shifts in temporal dynamics, evolving trends, and long-range and short effects
- Self-supervised approach: TF-C uses contrastive learning to inject TF-C into a pre-training model, bringing the encoded time-based and frequency-based representations and their local neighborhoods close together in the latent space
- Strong **generalization**: Our findings have implications for building broadly generalizable pre-training models for time series

Sequence-to-sequence neural models: Beyond natural language understanding

- State-of-the-art for sequence modeling
- Self attention
- No-recurrent units, allowing parallel computation
- Widely used in almost all language tasks now
 - Machine translation
 - Text generation
 - Question answering



Transformer architectures for time series

Autoformer – decomposing the time series components



Prediction

Transformer architectures for time series



arXiv:1912.09363

Temporal fusion transformer

Transformer architectures for time series



TSMixer – An all MLP architecture for time series forecasting

Today's lecture

Time series tasks in healthcare

2 Self-supervised pre-training for time series

3 Learning representations of regular and irregular time series

4 Understanding time series models

Irregular vs. regular time series

Regular time series



Irregular time series



Why are irregular time series challenging?

Prevailing methods:

- Assume aligned measurements
- Assume fixed-sized input data
- Impute or fill-in missing values

Irregular time series:

- Observations across sensors are not aligned
- Varying times among adjacent observations
- Arbitrary length: different samples have varyin
- Different subsets of sensors recorded



Problem definition

Input:

- Dataset D of irregular time series samples
- Every sample S_i can have multiple sensors
- Every sensor can have arbitrary number of irregularly sampled observations/readouts



- Raindrop learns a function $f: S_i \rightarrow z_i$ that maps S_i to a fixed-length representation z_i suitable for downstream tasks of interest, such as classification
- Using learned z_i , one can predict label $\hat{y}_i \in \{1, ..., C\}$ for S_i

Model irregularity by leveraging sensor dependencies

- Sensors are not independent of each other:
 - Inter-sensor dependencies contain useful information about time series
- Idea: Leverage relational structure among sensors
 - Learn latent graph structures from multivariate time series and model timevarying inter-sensor dependencies through neural message passing
 - Model sample-varying and time-varying relational structure in irregular time series

Next: What motivates the use of inter-sensor dependencies?

Raindrop: Observations as "raindrops" hitting a "surface"

- Observations (raindrops) hit the sensor graph (surface) asynchronously and at irregular times
- Observations are processed by passing messages to neighboring sensors (creating ripples), taking into account learned sensor dependencies





$$x_u = f(x_{v1}, x_{v2}, x_{v3})$$

Generate embedding of node *u* by capturing node dependencies through message passing

Sensor dependency graphs



Passing messages between neighboring sensors in every sample



Passing messages between neighboring sensors in every sample



Overview of Raindrop

Hierarchical learning

- Step 1: Construct sensor dependency graphs
 - For every sample, initialize a fully-connected graph
 - During training, update neighbors & edge weights:
 - Graphs are time-sensitive
 - Graphs are sample-sensitive
 - Similar graphs for similar samples
- Step 2: Sensor u is activated = its value is observed
 - Use message passing to generate observations for neighbors of active sensor u
- Step 3: <u>Sensor</u> embeddings
 - For sensor *u*, aggregate observation embeddings across all timestamps into *u*'s embedding
- Step 4: <u>Sample</u> embedding
 - Gather embeddings across all sensors into a representation of sample S_i using a readout function



Datasets & evaluation setup

P19: PhysioNet Sepsis Early Prediction

- 40,336 patients, 34 sensors
- Classification: Sepsis occurring or not
- P12: PhysioNet Mortality Prediction
 - 11,988 patients, 36 sensors
 - Classification: Length of stay in the ICU (>3 days or not)

- PAM: PAMAP2 Physical Activity Monitoring
 - 5,333 samples, 17 sensors
 - 8-class classification: 8 activities of daily lives







Setting 1/3: Time series classification

Given irregular sensor readouts of a given sample, predict a label for it

	P19		P12		PAM			
Methods	AUROC	AUPRC	AUROC	AUPRC	Accuracy	Precision	Recall	F1 score
Transformer	83.2 ± 1.3	47.6 ± 3.8	65.1 ± 5.6	95.7 ± 1.6	83.5 ± 1.5	84.8 ± 1.5	86.0 ± 1.2	85.0 ± 1.3
Trans-mean	84.1 ± 1.7	47.4 ± 1.4	66.8 ± 4.2	95.9 ± 1.1	83.7 ± 2.3	84.9 ± 2.6	86.4 ± 2.1	85.1 ± 2.4
GRU-D	83.9 ± 1.7	46.9 ± 2.1	67.2 ± 3.6	95.9 ± 2.1	83.3 ± 1.6	84.6 ± 1.2	85.2 ± 1.6	84.8 ± 1.2
SeFT	78.7 ± 2.4	31.1 ± 2.8	66.8 ± 0.8	96.2 ± 0.2	67.1 ± 2.2	70.0 ± 2.4	68.2 ± 1.5	68.5 ± 1.8
mTAND	80.4 ± 1.3	32.4 ± 1.8	65.3 ± 1.7	96.5 ± 1.2	74.6 ± 4.3	74.3 ± 4.0	79.5 ± 2.8	76.8 ± 3.4
IP-Net	84.6 ± 1.3	38.1 ± 3.7	72.5 ± 2.4	96.7 ± 0.3	74.3 ± 3.8	75.6 ± 2.1	77.9 ± 2.2	76.6 ± 2.8
DGM^2-O	86.7 ± 3.4	44.7 ± 11.7	71.2 ± 2.5	96.9 ± 0.4	82.4 ± 2.3	85.2 ± 1.2	83.9 ± 2.3	84.3 ± 1.8
MTGNN	81.9 ± 6.2	39.9 ± 8.9	67.5 ± 3.1	96.4 ± 0.7	83.4 ± 1.9	85.2 ± 1.7	86.1 ± 1.9	85.9 ± 2.4
RAINDROP	$\textbf{87.0} \pm \textbf{2.3}$	51.8 ± 5.5	72.1 ± 1.3	$\textbf{97.0} \pm \textbf{0.4}$	88.5 ± 1.5	$\textbf{89.9} \pm \textbf{1.5}$	$\textbf{89.9} \pm \textbf{0.6}$	89.8 ± 1.0

- Raindrop achieves strong performance across three benchmarks
- In binary classification (P19 and P12), it outperforms baselines by 5.3% in AUROC and 4.8% in AUPRC
- In a challenging 8-way classification (PAM), it outperforms baselines by 5.7% in accuracy and 5.5% in F1 score

Setting 2/3: Leave-sensors-out

Dataset P19:

- 38,803 patients, 34 sensors
- Label: Sepsis or not
- Missing sensors: 10-50%

Larger missing rate \rightarrow Large improvement in performance



Missing rate	Model	Accuracy	Precision	Recall	F1 score
	Transformer	60.3 ± 2.4	57.8 ± 9.3	59.8 ± 5.4	57.2 ± 8.0
	Trans-mean	60.4 ± 11.2	61.8 ± 14.9	60.2 ± 13.8	58.0 ± 15.2
10%	GRU-D	65.4 ± 1.7	72.6 ± 2.6	64.3 ± 5.3	63.6 ± 0.4
	SeFT	58.9 ± 2.3	62.5 ± 1.8	59.6 ± 2.6	59.6 ± 2.6
	mTAND	58.8 ± 2.7	59.5 ± 5.3	64.4 ± 2.9	61.8 ± 4.1
	RAINDROP	77.2 ± 2.1	$\textbf{82.3} \pm \textbf{1.1}$	$\textbf{78.4} \pm \textbf{1.9}$	75.2 ± 3.1
	Transformer	63.1 ± 7.6	71.1 ± 7.1	62.2 ± 8.2	63.2 ± 8.7
	Trans-mean	61.2 ± 3.0	74.2 ± 1.8	63.5 ± 4.4	64.1 ± 4.1
20%	GRU-D	64.6 ± 1.8	73.3 ± 3.6	63.5 ± 4.6	64.8 ± 3.6
	SeFT	35.7 ± 0.5	42.1 ± 4.8	38.1 ± 1.3	35.0 ± 2.2
	mTAND	33.2 ± 5.0	36.9 ± 3.7	37.7 ± 3.7	37.3 ± 3.4
	RAINDROP	66.5 ± 4.0	72.0 ± 3.9	67.9 ± 5.8	65.1 ± 7.0
	Transformer	31.6 ± 10.0	26.4 ± 9.7	24.0 ± 10.0	19.0 ± 12.8
	Trans-mean	42.5 ± 8.6	45.3 ± 9.6	37.0 ± 7.9	33.9 ± 8.2
30%	GRU-D	45.1 ± 2.9	51.7 ± 6.2	42.1 ± 6.6	47.2 ± 3.9
	SeFT	32.7 ± 2.3	27.9 ± 2.4	34.5 ± 3.0	28.0 ± 1.4
	mTAND	27.5 ± 4.5	31.2 ± 7.3	30.6 ± 4.0	30.8 ± 5.6
	RAINDROP	52.4 ± 2.8	60.9 ± 3.8	51.3 ± 7.1	$\textbf{48.4} \pm \textbf{1.8}$
	Transformer	23.0 ± 3.5	7.4 ± 6.0	14.5 ± 2.6	6.9 ± 2.6
	Trans-mean	25.7 ± 2.5	9.1 ± 2.3	18.5 ± 1.4	9.9 ± 1.1
40%	GRU-D	46.4 ± 2.5	64.5 ± 6.8	42.6 ± 7.4	44.3 ± 7.9
	SeFT	26.3 ± 0.9	29.9 ± 4.5	27.3 ± 1.6	22.3 ± 1.9
	mTAND	19.4 ± 4.5	15.1 ± 4.4	20.2 ± 3.8	17.0 ± 3.4
	RAINDROP	52.5 ± 3.7	53.4 ± 5.6	$\textbf{48.6} \pm \textbf{1.9}$	44.7 ± 3.4
	Transformer	21.4 ± 1.8	2.7 ± 0.2	12.5 ± 0.4	4.4 ± 0.3
	Trans-mean	21.3 ± 1.6	2.8 ± 0.4	12.5 ± 0.7	4.6 ± 0.2
50%	GRU-D	37.3 ± 2.7	29.6 ± 5.9	32.8 ± 4.6	26.6 ± 5.9
	SeFT	24.7 ± 1.7	15.9 ± 2.7	25.3 ± 2.6	18.2 ± 2.4
	mTAND	16.9 ± 3.1	12.6 ± 5.5	17.0 ± 1.6	13.9 ± 4.0
	RAINDROP	46.6 ± 2.6	44.5 ± 2.6	42.4 ± 3.9	38.0 ± 4.0

Setting 3/3: Group-wise classification

- Split by age: Patients older than 65 years vs. younger patients
- Split by gender: Male vs. female patients
- Use one group as training set and randomly split the other group into validation (50%) and test set (50%)

Generalizing to a new patient group

Model	Generalizing to a new patient group							
	Train: Young \rightarrow Test: Old Train: Old \rightarrow Test: Young			$ $ Train: Male \rightarrow Test: Female $ $ Train: Female \rightarrow Test: Male				
	AUROC	AUPRC	AUROC	AUPRC	AUROC	AUPRC	AUROC	AUPRC
Transformer Trans-mean GRU-D SeFT mTAND	$\begin{array}{c} 76.2 \pm 0.7 \\ 80.6 \pm 1.4 \\ 76.5 \pm 1.7 \\ 77.5 \pm 0.7 \\ 79.0 \pm 0.8 \end{array}$	$\begin{array}{c} 30.5 \pm 4.8 \\ 39.8 \pm 4.2 \\ 29.5 \pm 2.3 \\ 26.6 \pm 1.2 \\ 28.8 \pm 2.3 \end{array}$	$\begin{array}{c} 76.5 \pm 1.1 \\ 78.4 \pm 1.1 \\ 79.6 \pm 1.7 \\ 78.9 \pm 1.0 \\ 79.4 \pm 0.6 \end{array}$	$\begin{array}{c} 33.7 \pm 5.7 \\ 35.8 \pm 2.9 \\ 35.2 \pm 4.6 \\ 32.7 \pm 2.7 \\ 29.8 \pm 1.2 \end{array}$	$77.8 \pm 1.1 \\ 80.2 \pm 1.7 \\ 78.5 \pm 1.6 \\ 78.6 \pm 0.6 \\ 78.0 \pm 0.9$	$\begin{array}{c} 26.0 \pm 6.2 \\ 32.1 \pm 1.9 \\ 31.9 \pm 4.8 \\ 31.1 \pm 1.2 \\ 26.5 \pm 1.7 \end{array}$	$ \begin{array}{c c} 75.2 \pm 1.0 \\ 76.4 \pm 0.8 \\ 76.3 \pm 2.5 \\ 76.9 \pm 0.5 \\ 78.9 \pm 1.2 \end{array} $	30.3 ± 5.5 32.5 ± 3.3 31.1 ± 2.6 26.4 ± 1.1 29.2 ± 2.0
RAINDROP	$\textbf{83.2} \pm \textbf{1.6}$	$\textbf{43.6} \pm \textbf{4.7}$	$\textbf{82.0} \pm \textbf{4.4}$	$\textbf{44.3} \pm \textbf{3.6}$	$\textbf{85.0} \pm \textbf{1.4}$	$\textbf{45.2} \pm \textbf{2.9}$	81.2 ± 3.8	$\textbf{40.7} \pm \textbf{2.9}$

- Strong results in across-group scenarios; e.g., Raindrop outperforms strongest baseline by 4.8% in AUROC and 13.1% in AUPRC when asked to generalize to female patients
- Raindrop can generalize to new samples unseen during training:
 - Sensor dependency graphs are sample-specific and estimated using a sample's observations
 - Raindrop can adaptively generate sensor dependencies based on readouts of a test sample

Can Raindrop learn dynamics of sensors?

Nodes: sensors



Nodes 1 (pulse oximetry), 5 (diastolic BP), and 12 (partial pressure of carbon dioxide from arterial blood) have lower weights in patients with no sepsis

> Patients with no sepsis Sensor dependency graphs averaged across negative samples



Patients who develop sepsis Sensor dependency graphs averaged across positive samples

Distinguishable patterns between graphs of negative and positive samples → Raindrop learns dynamics of sensors purely from observational data

Can Raindrop learn dynamics of sensors?



Show is differential inter-sensor graph between patients who will likely develop sepsis and those who won't

Edges are colored by the divergences; darker colors denote sensor dependencies that are more crucial to patient classification

Irregular time series: Recap

- Irregular time series: Raindrop addresses the complexity of time series, e.g., misaligned observations, varying time gaps & varying numbers of observations per sensor
- Inter-sensor structure: Raindrop adopts neural message passing to model inter-sensor dependencies in irregular time series
- Great generalization: Raindrop has excellent performance in challenging settings, including setups where a subset of sensors malfunction (i.e., have no readouts at all)

Today's lecture

Time series tasks in healthcare

2 Self-supervised pre-training for time series

3 Learning representations of regular and irregular time series

Understanding time series models

Intrinsic vs. post-hoc interpretability

Intrinsic (e.g. linear models, trees, attention)



Mihaela van der Schaar

Standard feature importance scoring methods

Highlight most important features for the model

• Integrated Gradient [Sundararajan et al. 2017]

$$a_i(f, \mathbf{x}) = \left(x_i - x_i^0\right) \times \int_0^1 \frac{\partial f[\mathbf{x}^0 + t(\mathbf{x} - \mathbf{x}^0)]}{\partial x_i} dt$$



• SHAP [Lundberg et al. 2017]

$$a_i(f, \mathbf{x}) = \sum_{S \subset [\dim \mathcal{X}] \setminus \{i\}} \frac{|S|! (\dim \mathcal{X} - |S| - 1)}{(\dim \mathcal{X})!} [f(\mathbf{x}_{S \cup \{i\}}) - f(\mathbf{x}_S)]$$



"Standard" feature importance methods perform poorly for time-series [Ismail et al., NeurIPS 2020]

Mihaela van der Schaar

Challenges for explaining time series

Not easily visually interpretable

- Noisy samples
- Dense informative features, unlike imaging and text modalities

Temporal patterns

- Only show up when looking at time segments and long-term behaviors
- Perturbations matter
 - Setting a value to zero does not ignore that time point
 - Temporal dependencies cannot be ignored



What makes time series datasets different?



Mihaela van der Schaar

How to detect salient features?

Perturbation based detection

Premise: salient features affect the model's prediction

Detect salient features by feature perturbations

Feature perturbation affects prediction \rightarrow Salient feature



How to take the time context into account?

Time context matters

Typical saliency methods treat each input $x_{t,i}$ as a feature

 \Rightarrow Time dependency is ignored by the saliency method

Dynamic Perturbation Operator

Idea: perturb each $x_{t^*,i}$ by using neighbouring times:

 $\begin{array}{ll} \mbox{Perturbed input} & t^{*}+W_{2} \mbox{ Linear combination} \\ \pi(x_{t^{*},i}\,;t^{*},i) = & \sum_{t=t^{*}-W_{1}} c_{t}(t^{*},i) \times x_{t,i} \end{array}$

 \Rightarrow Time dependency is integrated in perturbation





How to take the time context into account?

Time context matters

Typical saliency methods treat each input $x_{t,i}$ as a feature

 \Rightarrow Time dependency is ignored by the saliency method

Dynamic Perturbation Operator

Idea: perturb each $x_{t^*,i}$ by using neighbouring times:

 $\begin{array}{ll} \mbox{Perturbed input} & t^{*}+W_{2} \mbox{ Linear combination} \\ \pi(x_{t^{*},i}\,;t^{*},i) = & \sum_{t=t^{*}-W_{1}} c_{t}(t^{*},i) \times x_{t,i} \end{array}$

 \Rightarrow Time dependency is integrated in perturbation





Dynamask



Backpropagate

How to make saliency masks parsimonious?

What do we mean by parsimonious?

Masks should not highlight more features than necessary

 \Rightarrow We need to enforce feature selection

How to enforce parsimony?

The user selects the fraction a of most important features

We add a regularization to enforce sparsity:

 $\mathcal{L}_{a}(\mathbf{M}) = \|vecsort(\mathbf{M}) - \mathbf{r}_{a}\|^{2}$

Sets the $(1-a) \times T \times d_X$ smallest mask coefficients to zero



How to avoid quick variation in the saliency mask?

Quick time variations of the saliency

Might want to avoid quick time variations of the saliency

This can be a prior belief or a preference of the user

How to avoid this?

We add a regularization to penalize saliency jumps over time:

$$\mathcal{L}_{c}(\mathbf{M}) = \sum_{t=1}^{T-1} \sum_{i=1}^{d_{X}} \left| m_{t+1,i} - m_{t,i} \right|$$



Dynamask - Example

Example number 5


Existing time series explainers are inadequate

Perturbations are continuous

- Can deform shape of samples
- 2 Give only instance-based explanations
 - Cannot relate patterns across samples
- 3 Fail to match performance of generic explainers
 - Post-hoc methods suffer from a lack of faithfulness and stability



Desiderata for time series explanations

- Temporally connected and visually digestible
- Identify the <u>location</u> of predictive time series signals and underlying interpretable <u>patterns</u>
- Connect explanations across samples

TimeX is a time-series consistency explainer

- Surrogate model to mimic the behavior of a pretrained time series model
- TimeX makes inferences on masked samples
- Model behavior consistency
 - Enforces faithfulness at the level of the latent space
 - Learns a flexible latent space of explanations



Encoding Time-Series Explanations through Self-Supervised Model Behavior Consistency, NeurIPS 2023

TimeX learns highly-faithful explanations

Method	AUPRC	SeqComb-MV AUP	AUR	AUPRC	LowVar AUP	AUR
IG Dynamask WinIT CoRTX SGT + Grad	$ \begin{array}{c c} 0.3298 \pm 0.0015 \\ 0.3136 \pm 0.0019 \\ 0.2809 \pm 0.0018 \\ 0.3629 \pm 0.0021 \\ \underline{0.4893} \pm 0.0005 \end{array} $	$\begin{array}{c} \underline{0.7483} {\pm} 0.0027 \\ 0.5481 {\pm} 0.0053 \\ 0.7594 {\pm} 0.0024 \\ 0.5625 {\pm} 0.0006 \\ 0.4970 {\pm} 0.0005 \end{array}$	$\begin{array}{c} 0.2581 {\pm} 0.0028 \\ 0.1953 {\pm} 0.0025 \\ 0.2077 {\pm} 0.0021 \\ 0.3457 {\pm} 0.0017 \\ \textbf{0.4289} {\pm} 0.0018 \end{array}$	$\begin{array}{c c} \textbf{0.8691} {\pm} 0.0035 \\ 0.1391 {\pm} 0.0012 \\ 0.1667 {\pm} 0.0015 \\ \underline{0.4983} {\pm} 0.0014 \\ 0.3449 {\pm} 0.0010 \end{array}$	$\begin{array}{c} \underline{0.4827} {\pm 0.0029} \\ 0.1640 {\pm 0.0028} \\ 0.1140 {\pm 0.0022} \\ 0.3281 {\pm 0.0027} \\ 0.2133 {\pm 0.0029} \end{array}$	$\begin{array}{c} \underline{0.8165} \pm 0.0016 \\ 0.2106 \pm 0.0018 \\ 0.3842 \pm 0.0017 \\ 0.4711 \pm 0.0013 \\ 0.3528 \pm 0.0015 \end{array}$
ТімеХ	0.6878 ±0.0021	0.8326 ±0.0008	0.3872 ± 0.0015	0.8673 ±0.0033	0.5451 ±0.0028	0.9004 ±0.0024

Method	AUPRC	ECG AUP	AUR
IG	0.4182 ± 0.0014	0.5949 ±0.0023	$0.3204 {\pm} 0.0012$
Dynamask	0.3280 ± 0.0011	$0.5249 {\pm} 0.0030$	$0.1082 {\pm} 0.0080$
WinIT	0.3049 ± 0.0011	$0.4431 {\pm} 0.0026$	0.3474 ± 0.0011
CoRTX	0.3735 ± 0.0008	$0.4968 {\pm} 0.0021$	0.3031 ± 0.0009
SGT + Grad	0.3144 ± 0.0010	$0.4241 {\pm} 0.0024$	$0.2639 {\pm} 0.0013$
TIMEX	0.4721 ±0.0018	0.5663 ± 0.0025	0.4457 ±0.0018

Encoding Time-Series Explanations through Self-Supervised Model Behavior Consistency, NeurIPS 2023

Landmarks learn important patterns in ECG



Landmarks partition the latent space of explanations into interpretable temporal patterns

Encoding Time-Series Explanations through Self-Supervised Model Behavior Consistency, NeurIPS 2023

Explanations through self-supervised model behavior consistency: Recap

- TimeX is a SOTA time series
 explainer
- STEs learn discrete masks that represent explanations and prevent shortcut learning
- Model behavior consistency preserves faithfulness across latent spaces
- Landmarks find important common temporal patterns to increase interpretability



Today's lecture

Time series tasks in healthcare

- 2 Self-supervised pre-training for time series
- 3 Learning representations of regular and irregular time series

4 Understanding time series models

Thank you for an incredible semester!

Final project presentations: See you all tomorrow (Wed) from 9am-12pm

AIM 2	Q Search AIM 2 Canvas BMI 702 Harvard DBMI Zitnik Lab					
Home	Artificial Intelligence in Medicine II					
Syllabus	Harvard - BMIE 202 and BMI 702 Spring 2025					
Course Project						
Focused Tutorials	Advances in Al will have a broad and protound impact on science and medicine, ortering new approaches to transform medical research and practice. This course provides a comprehensive overview of cutting-edge Al paradigms, including self-supervised learning, generative models, and multimodal techniques that integrate diverse data types. Beyond foundational methods, the course dives into a range of real-world applications in natural language processing, medical image analysis, relational and structure understanding, and longitudinal patient data.					
Calendar						
L1 - NLP I						
L2 - NLP II						
L3 - Generative Al	FACULTY INSTRUCTOR					
L4 - Agentic Al	Marinka Zitnik					
L5 - Medical Imaging I	marinka@hms.harvard.edu					
L6 - Medical Imaging II	Office Hours: Mon, 12pm – 1pm, Countway 309					
L7 - Trustworthy Al						
L8 - Networks I						
L9 - Networks II	https://zitniklab.hms.harvard.edu/AIM2					
L10 - Molecular Al						
L11 - Multimodal Al						
L12 - Ethical & Legal						
112 - Timo Sorios & Sonsors						

Schedule

This is just the beginning You now have the tools to bridge AI and medicine to ask deeper questions, to build more powerful models, and to reimagine what is possible for human health

The future of AI in medicine will be shaped by those who think critically, work collaboratively, and innovate responsibly. Keep pushing the boundaries of what we know. Stay grounded in science, but never lose sight of imagination

The discoveries you make can change lives

I look forward to seeing the breakthroughs you will create