

Learning Individualized Treatment Rules With Estimated Translated Inverse Propensity Score

Zhiliang Wu^{1,2}, Yinchong Yang¹, Yunpu Ma^{1,2}, Yushan Liu^{1,2}, Rui Zhao^{1,2}, Michael Moor³, Volker Tresp^{1,2}

¹Siemens AG, ²LMU Munich, ³ETH Zurich

IEEE International Conference on Healthcare Informatics, 2020, Best Paper Award



IEEE ICHI 2020



SIEMENS

Agenda

- Motivation
- Methods
- Experiments
- Summary and Outlook

Motivation

Predictive Modelling vs. Individualized Treatment Rules

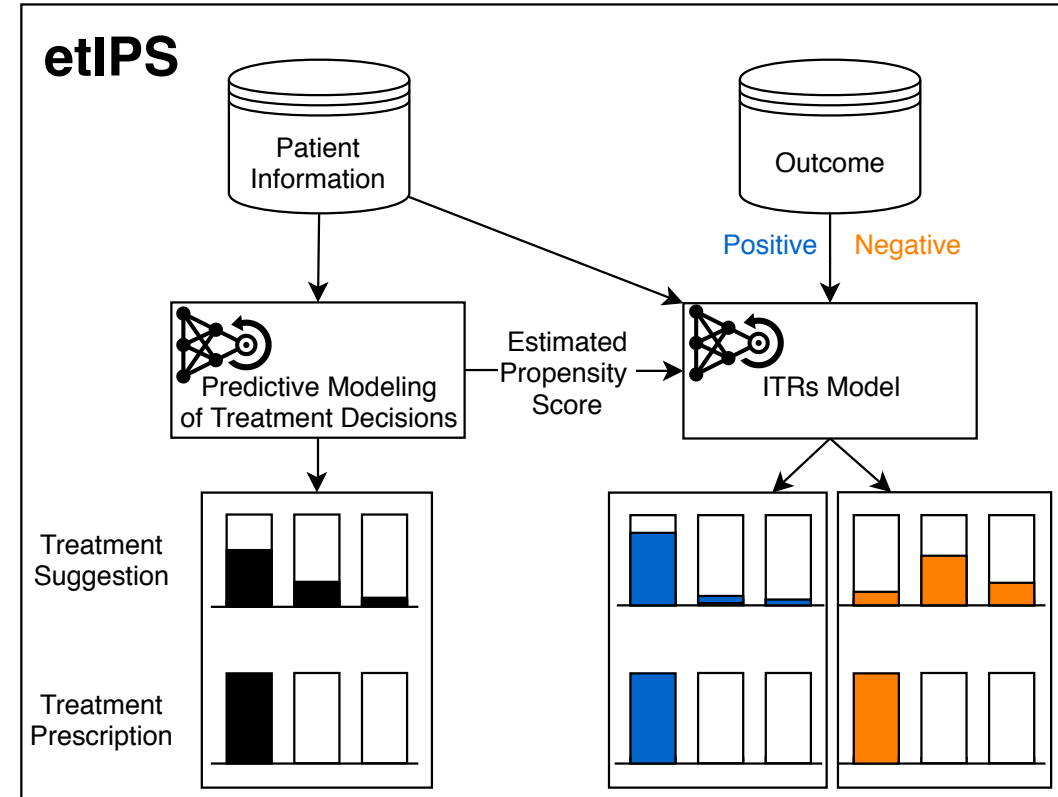
Predictive modelling of treatment decisions

[Esteban 2016, Yang 2017]

- Observed treatment decisions as targets
- Propensity score estimation
- No outcome information involved

Learning of the Individualized Treatment Rules (ITRs)

- A treatment policy that is expected to generate a better outcome for an individual patient
- Proposed framework: etIPS
- Integration of the observed outcome information with predictive modelling



Esteban, C., Staack, O., Baier, S., Yang, Y., & Tresp, V. (2016, October). Predicting clinical events by combining static and dynamic information using recurrent neural networks. In 2016 IEEE International Conference on Healthcare Informatics (ICHI) (pp. 93-101). IEEE.

Yang, Y., Fasching, P. A., & Tresp, V. (2017, August). Predictive modeling of therapy decisions in metastatic breast cancer with recurrent neural network encoder and multinomial hierarchical regression decoder. In 2017 IEEE International Conference on Healthcare Informatics (ICHI) (pp. 46-55). IEEE.

Motivation

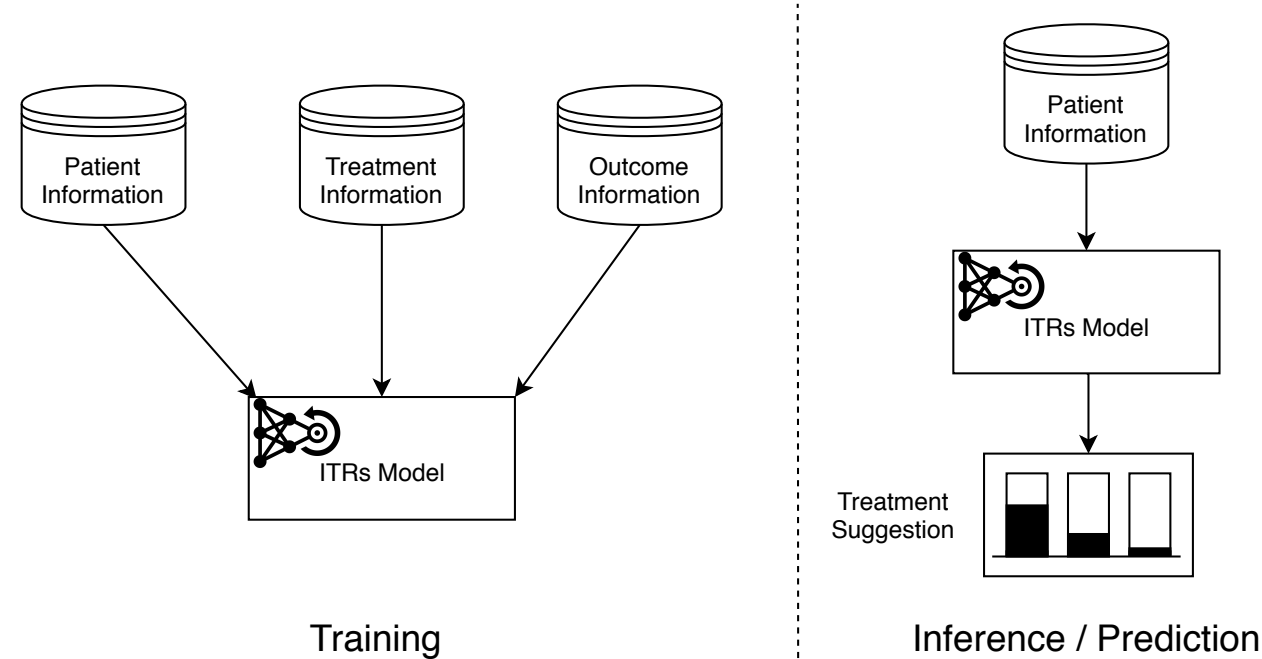
Problem Settings of ITRs

Training

- Treatment policy optimization
- Batch learning setup: No new samples collected
- Contextual bandits: Batch Learning from Bandit Feedback (BLBF)

Inference / Prediction

- The same as predictive modelling of treatment decisions



X_i : covariate of a patient
 a_i : treatment
 δ_i : loss

- Find a new policy π_w that minimizes the risk

$$\begin{aligned}
 r(\pi_w) &= \mathbb{E}_{X \sim \mathbb{P}(X)} \mathbb{E}_{a \sim \pi_w(a|X)} [\delta(X, a)] \\
 &= \mathbb{E}_{X \sim \mathbb{P}(X)} \mathbb{E}_{a \sim \mathbb{P}(a|X)} \left[\delta(X, a) \cdot \frac{\pi_w(a|X)}{\mathbb{P}(a|X)} \right]
 \end{aligned}$$

- Inverse Propensity Score Estimator (IPS)

$$\hat{r}_{\text{IPS}}(\pi_w) = \frac{1}{n} \sum_{i=1}^n \delta_i \frac{\pi_w(a_i | X_i)}{\mathbb{P}(a_i | X_i)}$$

- Self-Normalized IPS Estimator (SNIPS) [Swaminathan 2015]

$$\hat{r}_{\text{SNIPS}}(\pi_w) = \frac{\frac{1}{n} \sum_{i=1}^n \delta_i \frac{\pi_w(a_i | X_i)}{\mathbb{P}(a_i | X_i)}}{\frac{1}{n} \sum_{j=1}^n \frac{\pi_w(a_j | X_j)}{\mathbb{P}(a_j | X_j)}}$$

- λ -translated IPS estimator [Joachims2018]

$$\hat{r}_{\text{IPS}}^\lambda(\pi_w) = \frac{1}{n} \sum_{i=1}^n (\delta_i - \lambda) \frac{\pi_w(a_i | X_i)}{\mathbb{P}(a_i | X_i)}$$

Propensity score overfitting
[Swaminathan2015]

mini-batch
SGD friendly



The gradient of trainable parameters w in λ -translated IPS estimator

$$\hat{r}_{\text{IPS}}^\lambda(\pi_w) = \frac{1}{m} \sum_{i=1}^m (\delta_i - \lambda) \frac{\pi_w(a_i|X_i)}{\mathbb{P}(a_i|X_i)} \implies \nabla \hat{r}_{\text{IPS}}^\lambda = \frac{1}{m} \sum_{i=1}^m \frac{\delta_i - \lambda}{\mathbb{P}(a_i|X_i)} \nabla \pi_w(a_i|X_i)$$

For loss δ_i

- 0 for treatments with positive outcome
- 1 for treatments with negative outcome

Two extremes of λ -translated IPS estimator

- $\lambda = 0$: only the gradient of negative outcome is received by w . \rightarrow a completely different policy from physicians
- $\lambda = 1$: only the gradient of positive outcome is received by w . \rightarrow a very similar policy to the physicians
- $\lambda \in (0, 1)$: gradient of both negative and positive outcome can be received by w .
- The optimal λ^* : found through grid search.

Algorithm 1: etIPS

Input: A dataset of the form $\{X_i, a_i, \delta_i\}_{i=1}^m$.

Output: The policy of the optimal ITRs $\pi_{\mathbf{w}^*}(a|X)$.

1 Learn the physicians' policy $\hat{\mathbb{P}}(a|X)$ with $\{X_i, a_i\}_{i=1}^m$ using the network structure in Fig. 3.

2 Compute the estimated propensity score $\hat{p}_i := \hat{\mathbb{P}}(a = a_i|X_i)$ for all i .

3 **for** $\lambda_j \in (0, 1)$ **do**

$$4 \quad \left| \quad \mathbf{w}_j^* \leftarrow \arg \min_{\mathbf{w}_j} \left\{ \frac{1}{m} \sum_{i=1}^m (\delta_i - \lambda_j) \frac{\pi_{\mathbf{w}_j}(a_i|X_i)}{\hat{p}_i} \right\} \right.$$

$$5 \quad \left| \quad s_j \leftarrow \frac{1}{m} \sum_{i=1}^m \frac{\pi_{\mathbf{w}_j^*}(a_i|X_i)}{\hat{p}_i} \right.$$

6 **end**

$$7 \quad s^*, \mathbf{w}^* \leftarrow \arg \min_{s_j, \mathbf{w}_j^*} \left\{ \frac{1}{s_j} \frac{1}{m} \sum_{i=1}^m \delta_i \frac{\pi_{\mathbf{w}_j^*}(a_i|X_i)}{\hat{p}_i} \right\}$$

8 **return** $\pi_{\mathbf{w}^*}(a|X)$

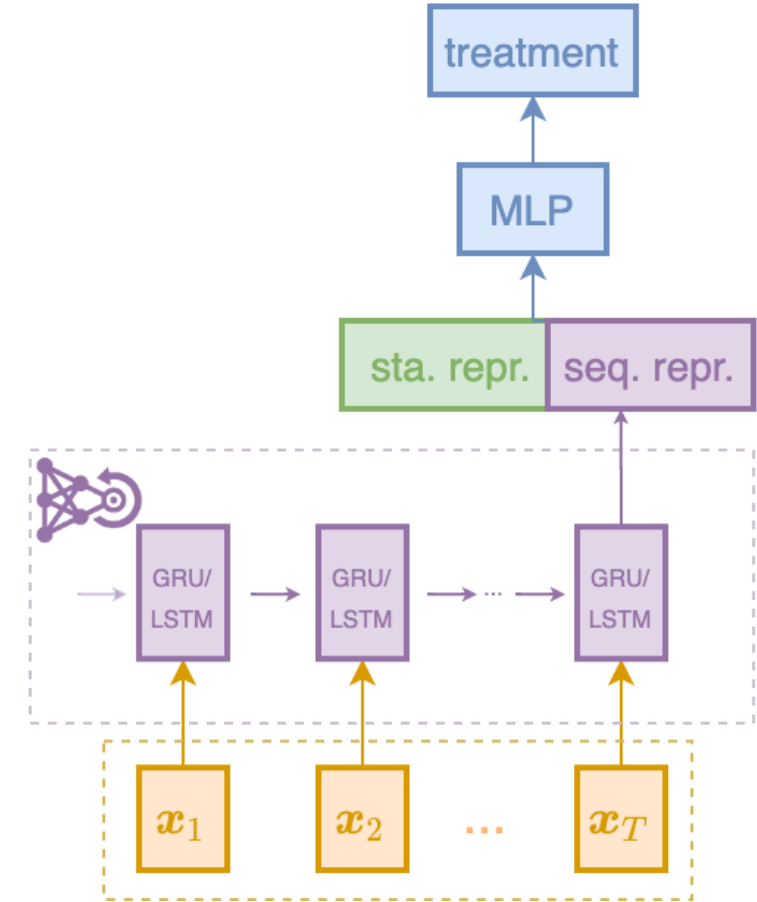
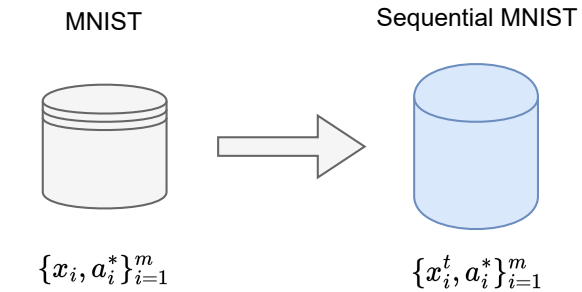


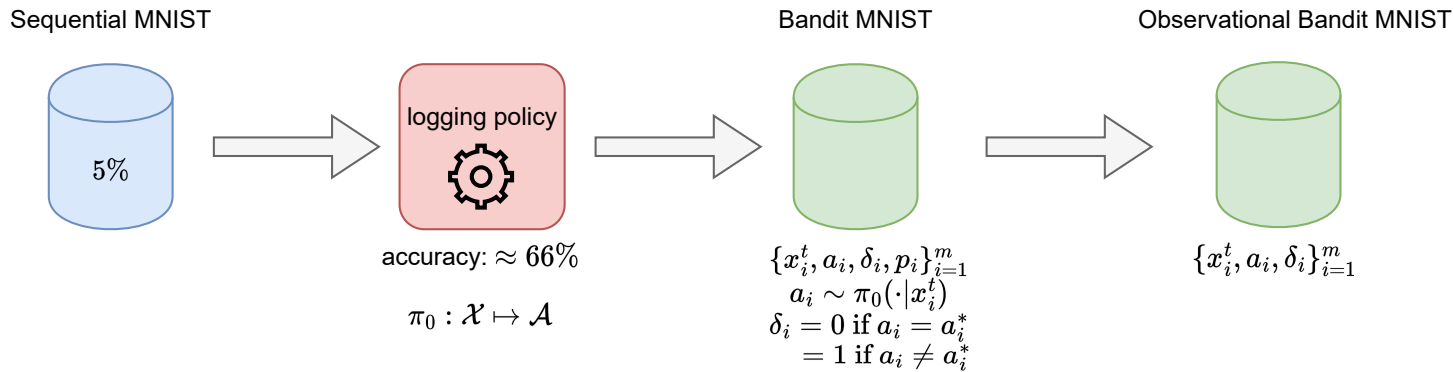
Fig. 3 Network architecture, first proposed in [Yang 2017]

Experiments

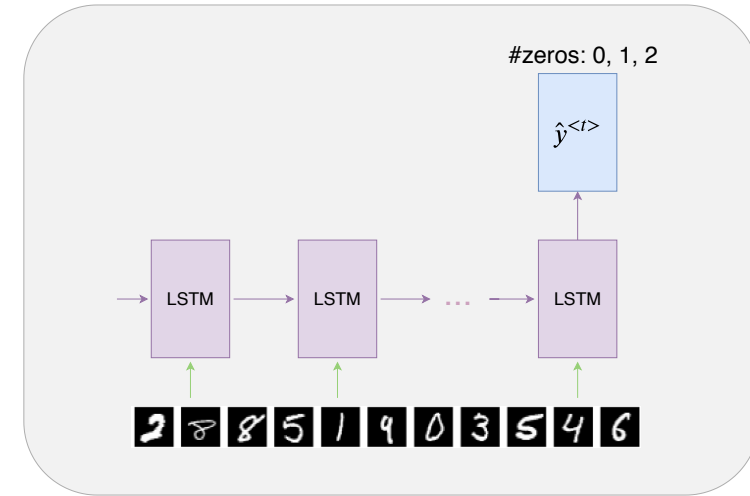
Simulation Study



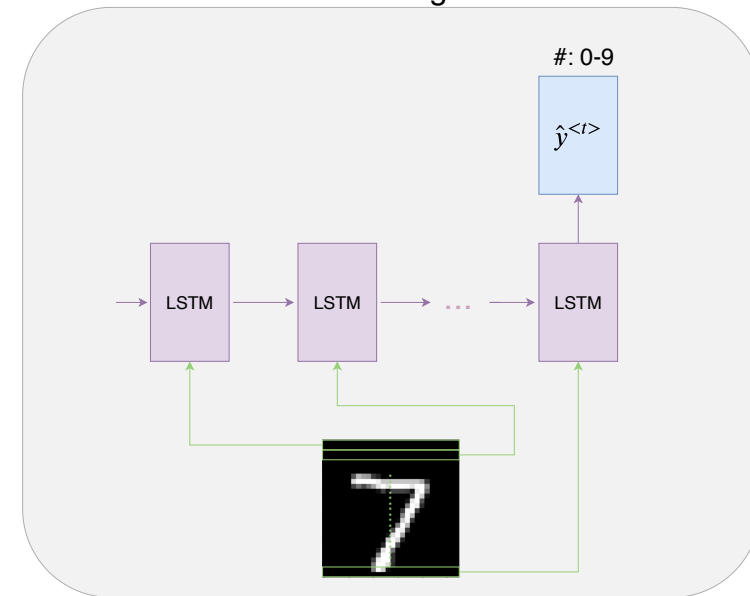
- Context: A sequence of images
- Action: Label prediction of the given sequence
- Loss: Correctness of the prediction
- Propensity: Probability of the respective label prediction



Supervised to Bandit Conversion Method [Agarwal 2014]



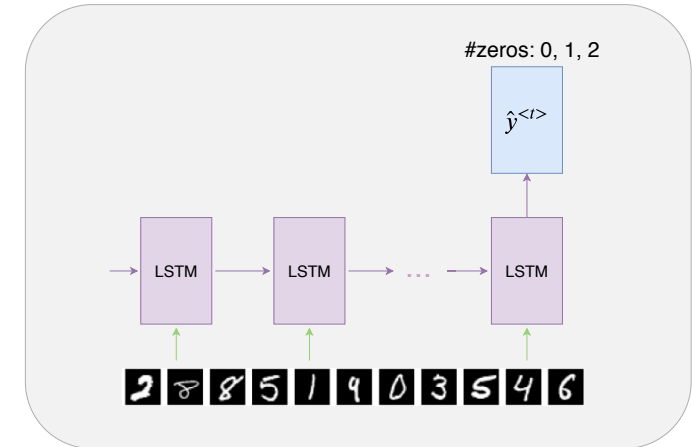
zero counting MNIST



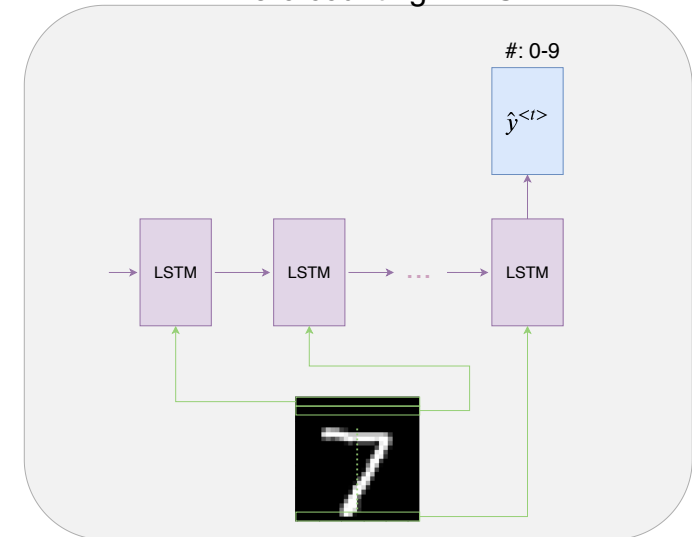
row-by-row MNIST

Hypothesis validation

- Without ground truth labels, the propensity score-based objective function is applicable to classification tasks with sequential data.
- The estimated propensity score could be used to replace the true propensity score in the objective function. (cf. *MLIPS* in [Xie 2018])



zero counting MNIST



row-by-row MNIST

TABLE I
ACCURACY OF DIFFERENT APPROACHES ON SEQUENTIAL CLASSIFICATION TASKS

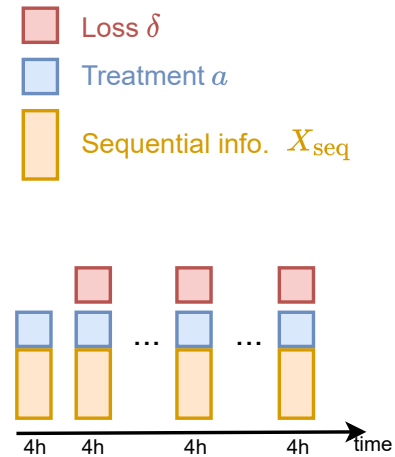
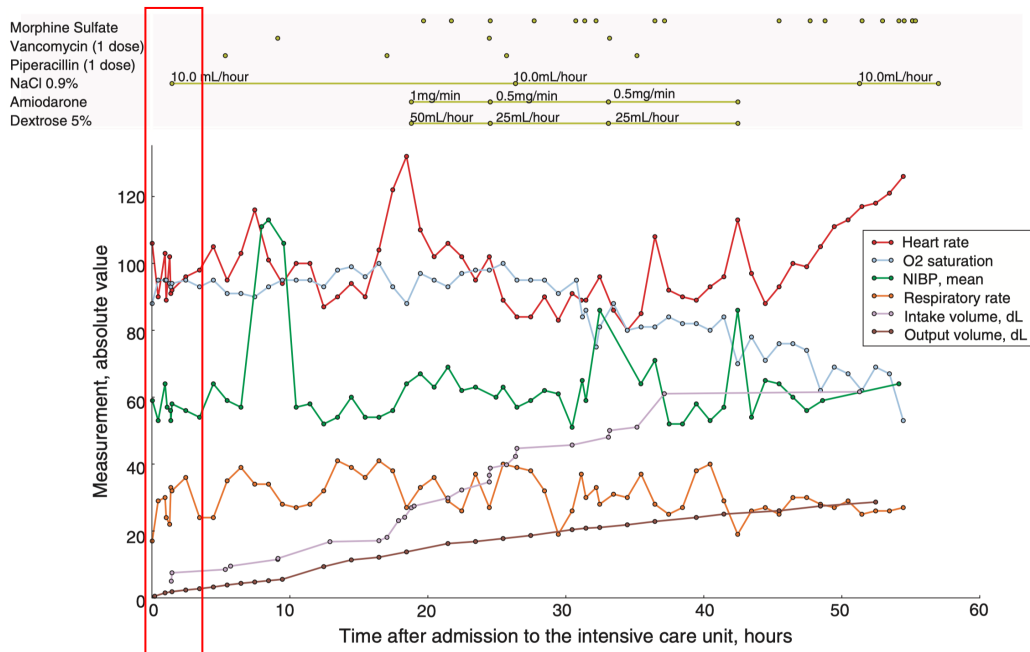
	propensity score	zeros counting MNIST	row-by-row MNIST
DM	-*	0.343 ± 0.0001	0.098 ± 0.0022
RP	-*	0.363 ± 0.0001	0.103 ± 0.0001
IPS	true	0.301 ± 0.012	0.020 ± 0.0061
tIPS	true	0.899 ± 0.0229	0.931 ± 0.0852
eIPS	estimated	0.319 ± 0.0075	0.016 ± 0.0098
etIPS	estimated	0.923 ± 0.0122	0.953 ± 0.0390

*The propensity score is not involved in the algorithm.

Medical Information Mart for Intensive Care database (MIMIC-III)

- A freely accessible database
- 53,423 Intensive Care Unit (ICU) admissions

Use case from [Komorowski 2018]



Dosage actions

		Dose of vasopressor				
		1	2	3	4	5
Dose of i. v. fluid	1	1	2	3	4	5
	2	6	7	8	9	10
	3	11	12	13	14	15
	4	16	17	18	19	20
	5	21	22	23	24	25

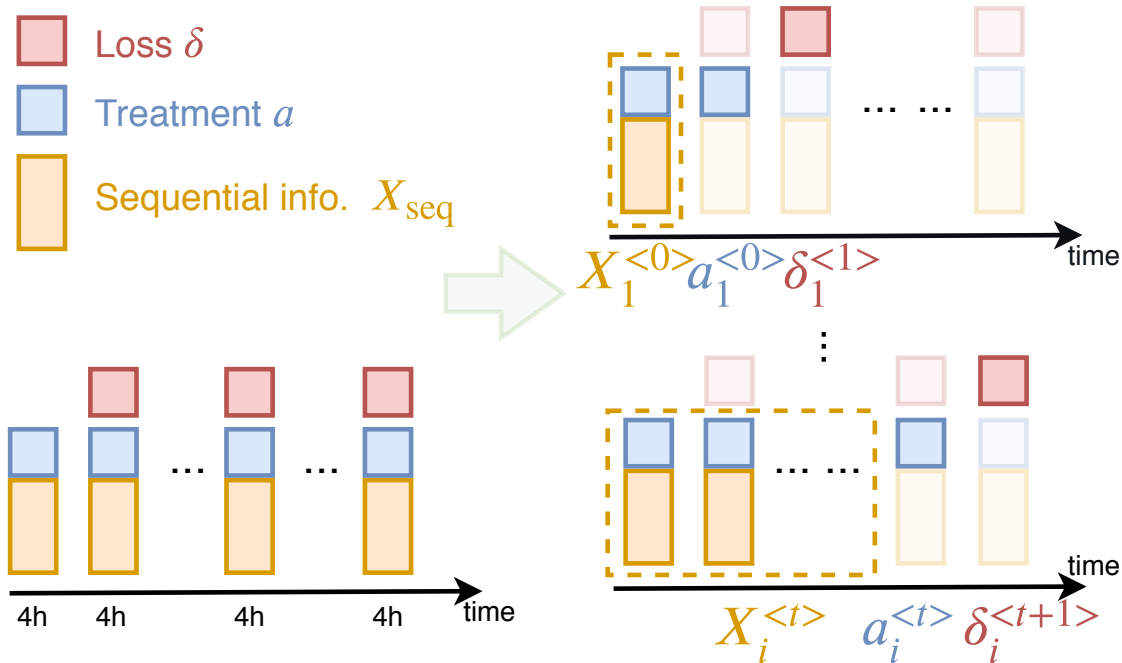


Experiments

MIMIC-III

Use case from [Komorowski 2018]

- A cohort of patients fulfilling the Sepsis-3 criteria [1]
- 20,944 admissions
- 224,333 samples



Different evaluation metrics to indicate the expected risk (lower is better)

- Average Treatment Effects under the new policy (ATENP) [Zhou 2017, Komorowski 2018]
- Doubly Robust Estimator (DR) [Dudik 2011]

TABLE II
EVALUATION WITH DIFFERENT RISK ESTIMATORS

	ATENP	DR
Predictive Modeling	-0.019 ± 0.0021	0.523 ± 0.0021
Direct Method*	0.032 ± 0.0001	-
Most Frequent [†]	-0.023 ± 0.0001	-
Random Policy	-0.023 ± 0.0001	0.478 ± 0.0026
Estimated Inverse Propensity Score	-0.025 ± 0.1009	0.504 ± 0.0071
Estimated Translated Inverse Propensity Score	-0.143 ± 0.0099	0.471 ± 0.0060

*There is no probability given by $\arg \min_a \mathbb{E}[\delta|X, a]$. The values for DR can therefore not be computed.

[†]A deterministic policy to suggest the most frequent treatment. No probability information is involved.

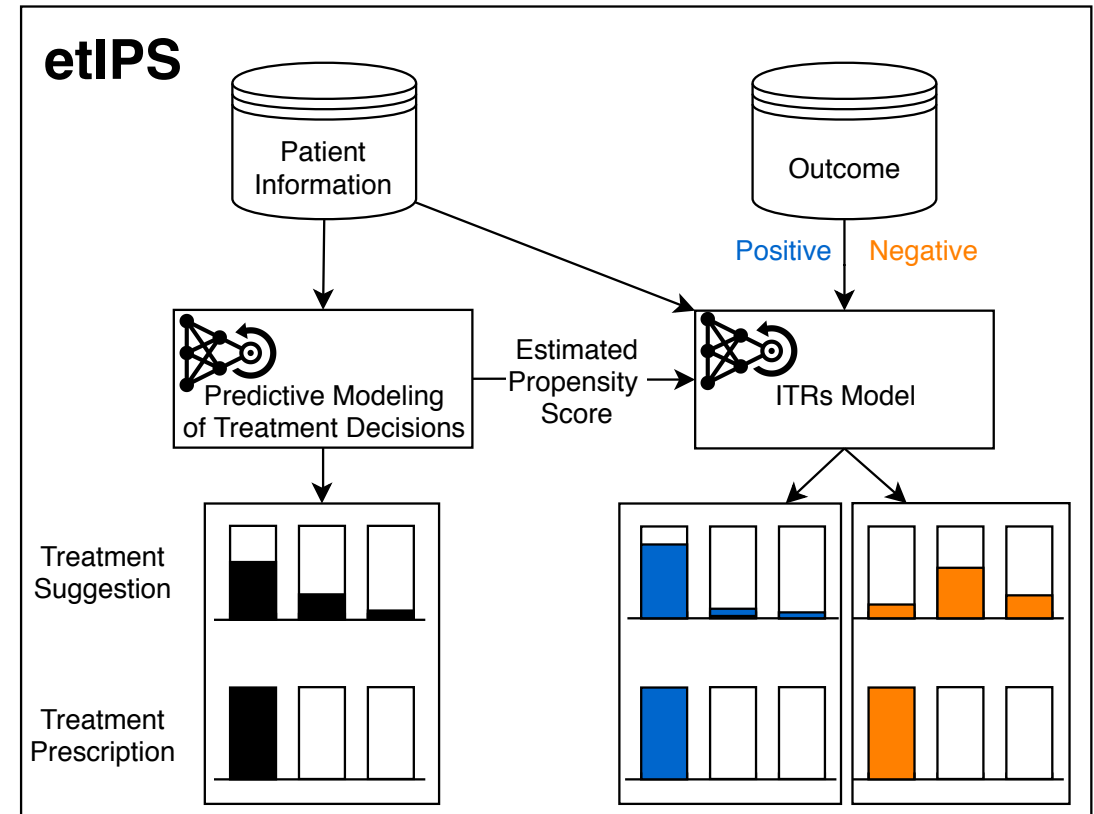
Komorowski, M., Celi, L. A., Badawi, O., Gordon, A. C., & Faisal, A. A. (2018). The artificial intelligence clinician learns optimal treatment strategies for sepsis in intensive care. *Nature medicine*, 24(11), 1716-1720.

Zhou, X., Mayer-Hamblett, N., Khan, U., & Kosorok, M. R. (2017). Residual weighted learning for estimating individualized treatment rules. *Journal of the American Statistical Association*, 112(517), 169-187.

Dudik, M., Langford, J., & Li, L. (2011, June). Doubly robust policy evaluation and learning. In *Proceedings of the 28th International Conference on International Conference on Machine Learning* (pp. 1097-1104).

Summary and Outlook

- A flexible framework to learn optimal ITRs
- Batch learning setup in healthcare domain
- Propensity score estimation
 - Predictive modelling of treatment decisions
- ITRs learning
 - Batch learning from bandit feedback (BLBF) algorithms
- Integration with potential outcome models
- Interpretability / Explanability of the learnt ITRs



| Thank you!

Dr. rer. nat. Zhiliang Wu
Research Scientist

Siemens AG
Otto-Hahn-Ring 6
81739 München

E-mail zhiliang.wu@siemens.com