# Evaluating RL Policies in Healthcare

**Online Evaluation**

High-stakes environment
- Potentially unsafe to patients
- Disruptive to human users
  and clinical workflows

Wiens et al. "Do no harm: a roadmap for responsible machine learning for health care." *Nature Medicine* 2019.

**Counterfactual-Augmented Importance Sampling for Semi-Offline Policy Evaluation**
Shengpu Tang, Jenna Wiens. NeurIPS 2023.

2

# Evaluating RL Policies in Healthcare

**Offline Evaluation**

**Online Evaluation**

?

Observational dataset
- Limited by available data
- May not reflect distribution shift induced by new policies
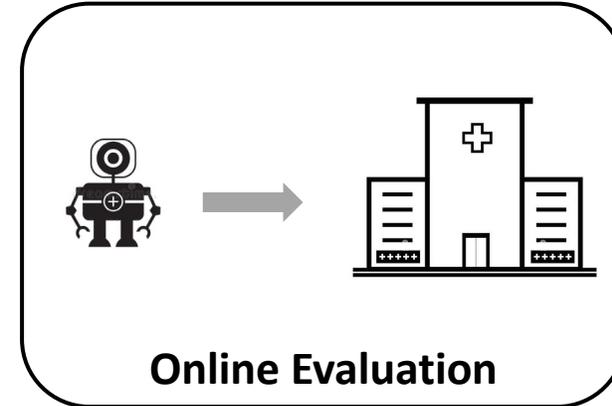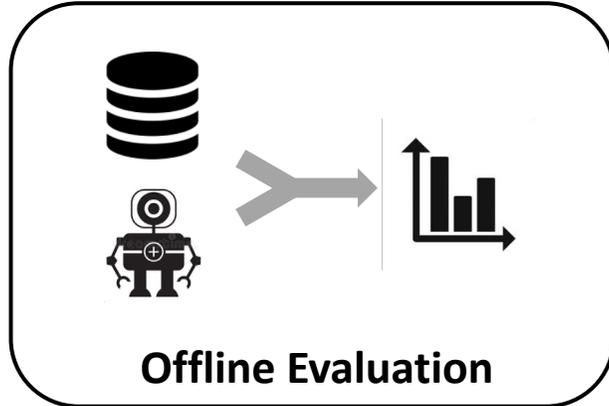
High-stakes environment
- Potentially unsafe to patients
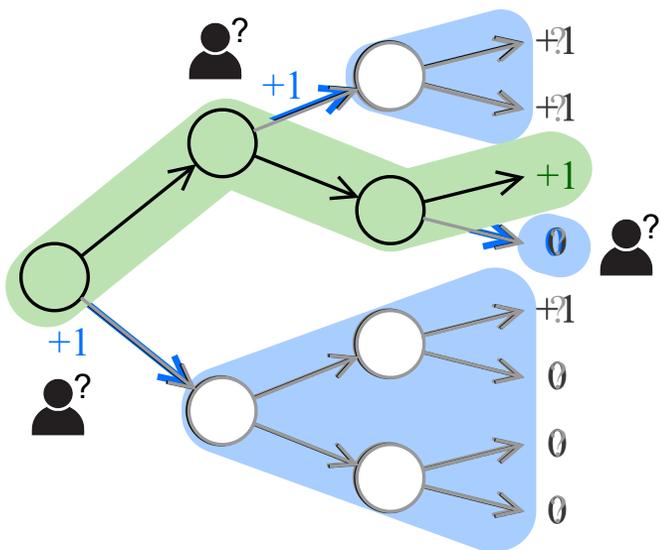- Disruptive to human users and clinical workflows

Wiens et al. "Do no harm: a roadmap for responsible machine learning for health care." *Nature Medicine* 2019.
Gottesman et al. "Guidelines for reinforcement learning in healthcare." *Nature Medicine* 2019.

**Counterfactual-Augmented Importance Sampling for Semi-Offline Policy Evaluation**
Shengpu Tang, Jenna Wiens. NeurIPS 2023.

3

# Our Contributions
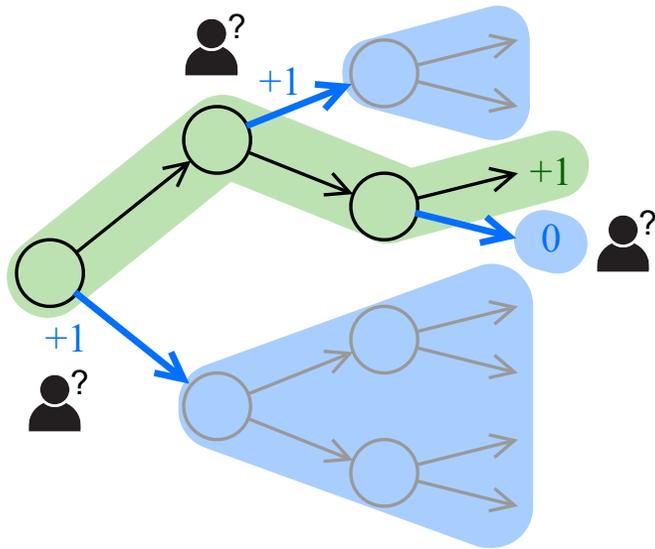
We propose a **semi-offline evaluation scheme** that combines observational data with **human annotations** of counterfactuals



Observational data contains **factual trajectories**

Query domain experts for **annotations**

of the **counterfactual trajectories**

**Counterfactual-Augmented Importance Sampling for Semi-Offline Policy Evaluation**
Shengpu Tang, Jenna Wiens. NeurIPS 2023.

4

# Augmenting Factual Data with Counterfactuals



Intuition: as if we collected **more data**

How do we use both **counterfactual annotations** and **observational data** to evaluate policies?

*"Simply adding annotations as new data"*
... is not theoretically valid.

**Counterfactual-Augmented Importance Sampling for Semi-Offline Policy Evaluation**
Shengpu Tang, Jenna Wiens. NeurIPS 2023.

5

where $\rho = \frac{\pi_e(a|s)}{\pi_b(a|s)}$



$$\hat{v}^{\mathrm{IS}} = \rho\, r$$

**Counterfactual-Augmented Importance Sampling for Semi-Offline Policy Evaluation**
Shengpu Tang, Jenna Wiens. NeurIPS 2023.

6

# Key Idea: Reweighted IS with Counterfactuals

where $\rho^{\tilde{a}} = \frac{\pi_e(\tilde{a}|s)}{\pi_b(\tilde{a}|s)}$

$$\hat{v}^{\text{C-IS}} = w^a \underbrace{\rho^a r}_{} + \sum_{\tilde{a} \in \mathcal{A} \setminus \{a\}} w^{\tilde{a}} \underbrace{\rho^{\tilde{a}} g^{\tilde{a}}}_{}$$

factual IS estimate
from observed reward

counterfactual IS estimates
based on annotations

where $w^a + \sum_{\tilde{a} \in \mathcal{A} \setminus \{a\}} w^{\tilde{a}} = 1$

**Counterfactual-Augmented Importance Sampling for Semi-Offline Policy Evaluation**
Shengpu Tang, Jenna Wiens. NeurIPS 2023.

7

# Theoretical Insights

$$\hat{v}^{\text{C-IS}} = w^a \rho^a r + \sum_{\tilde{a} \in \mathcal{A} \setminus \{a\}} w^{\tilde{a}} \rho^{\tilde{a}} g^{\tilde{a}}$$

Intuition: as if we collected **more data**
- More data for regions that lack support $\rightarrow$ reduce bias
- Even more data for regions with support $\rightarrow$ reduce variance

C-IS can achieve **lower bias** and **lower variance** than IS

**Counterfactual-Augmented Importance Sampling for Semi-Offline Policy Evaluation**
Shengpu Tang, Jenna Wiens. NeurIPS 2023.

8

# Experimental Results

Experiments conducted on the sepsis simulator

Based on the sepsis simulator introduced by Oberst & Sontag, ICML 2019.

**Compare**
- Standard approach (PDIS)
- Proposed approach (C-PDIS)

**Simulate collection of**
- Factual dataset
- Counterfactual annotations

to evaluate multiple treatment policies.

**Metrics**
- ↓ Evaluation error (RMSE)
- ↑ Ranking ability (Spearman correlation)

with respect to ground-truth policy performance

**Counterfactual-Augmented Importance Sampling for Semi-Offline Policy Evaluation**
Shengpu Tang, Jenna Wiens. NeurIPS 2023.

9

# Experimental Results

| Estimator | ↓ Evaluation Error | ↑ Ranking Ability |
|-----------|:------------------:|:-----------------:|
| **Baseline** | 0.113 | 0.596 |
| **Proposed** | 0.013 | 0.995 |

Our proposed approach **outperforms** the baseline method (without annotations) in terms of all metrics.

(under the assumption that annotations are "good")

**Counterfactual-Augmented Importance Sampling for Semi-Offline Policy Evaluation**
Shengpu Tang, Jenna Wiens. NeurIPS 2023.

10

# Experimental Results

| Estimator | ↓ Evaluation Error | ↑ Ranking Ability |
|---|---|---|
| **Baseline** | 0.113 | 0.596 |
| **Proposed** | 0.013 | 0.995 |
| **Proposed (biased)** | 0.028 | 0.979 |
| **Proposed (noisy)** | 0.029 | 0.977 |
| **Proposed (missing)** | 0.067 | 0.823 |

Our proposed approach **remains competitive** to the baseline method even with imperfect annotations (biased, noisy, missing).

**Counterfactual-Augmented Importance Sampling for Semi-Offline Policy Evaluation**
Shengpu Tang, Jenna Wiens. NeurIPS 2023.

11

# Takeaways

We propose a **new estimator** for **semi-offline evaluation** that combines observational data with **human annotations** of counterfactuals

$$\hat{v}^{\text{C-IS}} = w^a \rho^a r + \sum_{\tilde{a} \in \mathcal{A} \setminus \{a\}} w^{\tilde{a}} \rho^{\tilde{a}} g^{\tilde{a}}$$

- Theoretical insights show potential to *reduce both bias and variance*
- Experiments demonstrate robustness to *bias, noise, and missingness* of annotations

🔗 https://github.com/MLD3/CounterfactualAnnot-SemiOPE

**Counterfactual-Augmented Importance Sampling for Semi-Offline Policy Evaluation**
Shengpu Tang, Jenna Wiens. NeurIPS 2023.

12