

# Learning "What-if" Explanations For Sequential Decision Making

IE 708 Project

Hanish Prashant Dhanwalkar

November 19, 2024

# Batch IRL

## What is IRL:

- Learning technique that aims to infer the reward function of an agent by observing its behavior in a given environment.
- Unlike traditional reinforcement learning, where the reward function is explicitly defined, Batch IRL learns the reward function from a set of expert demonstrations or trajectories.

## Why Batch IRL

- **Online Learning:** Classic IRL algorithms require interactive access to the environment, or full knowledge of the environment's dynamics.
- Limited by the assumption that state dynamics are fully-observable and Markovian. NOT true for domains like medicine, where treatment depends on how patient covariates (tumour, side effects) have evolved over time.

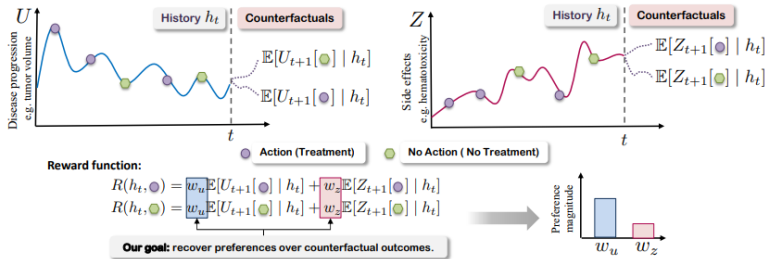
# "What-if" Explanations

Incorporating counterfactual reasoning into batch IRL

- Learn a parameterized reward function  $R(h_t, a_t)$  that is defined as a weighted sum over potential outcomes for taking action at given history  $h_t$ .
- This helps in reasoning out why an expert (eq. doctor) have chosen a particular action and "what" would have happen if any other alternative was taken. These experimentation are not possible in medicine domain.

# Example: Cancer Treatment

- Consider the decision making process of assigning a binary action given,
  - ▶ Tumour volume ( $U$ )
  - ▶ Side effects ( $Z$ )
- Let  $\mathbb{E}[U_{t+1}[a_t]|h_t]$  and  $\mathbb{E}[Z_{t+1}[a_t]|h_t]$  be the counterfactual outcomes for the two covariates when action  $a_t$  is taken given the history  $h_t$  of covariates and previous actions.
- The reward as the weighted sum of these counterfactuals:
  - ▶  $R(h_t, a_t) = w_u \mathbb{E}[U_{t+1}[a_t]|h_t] + w_z \mathbb{E}[Z_{t+1}[a_t]|h_t]$ ,
- This allows us to model the preferences of experts:  
e.g. finding that  $|w_u| > |w_z|$  indicates that the expert is treating more aggressively, by placing more weight on reducing tumour volume than on minimizing side effects.



**Figure:** Explaining decision-making behaviour in terms of preferences over “what if” outcomes. Evolution of tumour volume ( $U$ ) and side effects ( $Z$ ) under a binary action.

# Problem Formulation

At timestep  $t$ , let  $X_t \in X$  be the observed patient features and  $A_t \in A$  be the action (e.g. treatment) taken. Let  $x_t$  and  $a_t$  be realizations of these random variables,  $h_t = (x_0, a_0, \dots, x_{t-1}, a_{t-1}, x_t) = (x_{0:t}, a_{0:t-1}) \in H$  be a realization of the history  $H_t \in H$  of patient observations and actions.

- A policy  $\pi : H \times A \rightarrow [0, 1]$ , where  $\pi(a|h)$  indicates the probability of choosing action  $a \in A$  given history  $h \in H$  and  $\sum_{a \in A} \pi(a|h) = 1$ .
- Taking action  $a_t$  under history  $h_t$  results in observing  $x_{t+1}$  and obtaining  $h_{t+1}$ . The reward function is  $R : H \times A \rightarrow \mathbb{R}$  where  $R(h, a)$  represents the reward for taking action  $a \in A$  given history  $h \in H$ .
- The value function of a policy  $\pi$ ,  $V : H \rightarrow \mathbb{R}$  is defined as:  
 $V^\pi(h) = \mathbb{E}[\sum_{t=0}^{\infty} \gamma^t R(H_t, A_t) | \pi, H_0 = h]$ , where  $\gamma \in [0, 1]$  is the discount factor.
- The action-value function  $Q : H \times A \rightarrow \mathbb{R}$  of a policy is defined as  
 $Q^\pi(h, a) = \mathbb{E}[\sum_{t=0}^{\infty} \gamma^t R(H_t, A_t) | \pi, H_0 = h, A_0 = a]$

# Problem Formulation

- **Batch IRL:** consider a linear reward function  $R(h_t, a_t) = w \cdot \phi(h_t, a_t)$  where  $\|w\|_1 \leq 1$
- $\pi_E$  is attempting to optimise some unknown reward function  $R^*(h_t, a_t) = w^* \cdot \phi(h_t, a_t)$  where  $w^*$  are the 'true' reward weights.
- The value of policy  $\pi$  can be re-written as:

$$\mathbb{E}[V^\pi(H_0)] = w \cdot \mathbb{E}\left[\sum_{t=0}^{\infty} \gamma^t \phi(H_t, A_t) \mid \pi\right]$$

- The feature expectation of  $\pi$ , defined as the expected discounted cumulative feature vector obtained when choosing actions according to  $\pi$  is

$$\mu^\pi = \mathbb{E}\left[\sum_{t=0}^{\infty} \gamma^t \phi(H_t, A_t) \mid \pi\right] \in \mathbb{R}^d$$

such that:  $\mathbb{E}[V^\pi(H_0)] = w \cdot \mu^\pi$

# Problem Formulation

- Our aim is to recover the expert weights  $W^*$  as well as find a policy  $\pi$  that is close to the policy of the expert  $\pi_E$
- max-margin IRL approach and measure the similarity between the feature expectations of the expert's policy and the feature expectations of a candidate policy using  $\|\mu^{\pi_E} - \mu^{\pi}\|_2$
- In this batch IRL setting, we do not have knowledge of transition dynamics and we cannot sample more trajectories from the environment.



# Problem Formulation

- **Counterfactual reasoning:** To explain the expert's behaviour in terms of their trade-off associated with "what if" outcomes
- define the feature map  $\phi(h_t, a_t)$  part of the reward  
 $R(h_t, a_t) = w \cdot \phi(h_t, a_t)$
- Let  $Y[a]$  be potential outcome, either factual or counterfactual, for treatment  $a \in A$ . Using Dataset D, learn feature map  $\phi(h_t, a_t)$  such that  $\phi(h_t, a_t) = \mathbb{E}[Y_{t+1}[a_t] | h_t]$ .
- The potential outcomes for the other actions are the counterfactual ones and they allow us to understand what would happen to the patient if they receive a different treatment.
- Consider the model for estimating counterfactuals as a black box such that the feature map  $\phi$  represents the effect of taking action  $a_t$  for history  $h_t$ .

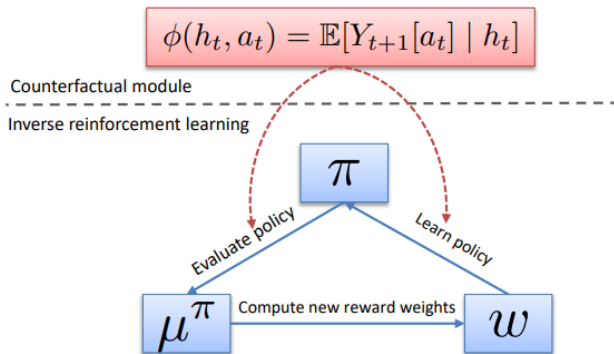
$$R(h_t, a_t) = w \cdot \phi(h_t, a_t) = w \cdot \mathbb{E}[Y_{t+1}[a_t] | h_t]$$

# Batch IRL using Counterfactuals

- Max-margin IRL starts with an initial random policy  $\pi$  and iteratively performs the following steps to recover the expert policy and its reward weights:
  - 1 estimate feature expectations  $\mu^\pi$  of candidate policy  $\pi$ ,
  - 2 compute new reward weights  $w$ ,
  - 3 find new candidate policy  $\pi$  that is optimal for reward function  $R$
- This approach finds a policy  $\tilde{\pi}$  that satisfies  $\|\mu^{\pi^E} - \mu^\pi\|_2 < \epsilon$  such that  $\tilde{\pi}$  has an expected value function close the expert policy.
- The expert feature expectations can be estimated empirically from the dataset  $D$  using:

$$\mu^{\pi^E} = \frac{1}{N} \sum_{i=0}^N \sum_{t=0}^{T^i} \gamma^t \phi(h_t^i, a_t^i)$$

# Batch IRL using Counterfactuals



**Figure:** CIRL. Counterfactuals are used to define  $\phi(h, a)$ , to estimate feature expectations  $\mu^\pi$  of candidate policy  $\pi$  in batch setting and to learn optimal policy for reward weights  $w$ .

# Counterfactual $\mu$ - Learning

First action  $a$  is taken randomly and for  $t \geq 1$ ,  $A_t \sim \pi(\cdot | H_t)$ .

This can be re-written as:

$$\begin{aligned}\mu^\pi(h, a) &= \phi(h, a) + \mathbb{E}_{h', a' \sim \pi(\cdot | h')} \left[ \sum_{t=0}^{\infty} \gamma^t \phi(h_t, a_t) | \pi, H_1 = h', A_1 = a' \right] \\ &= \phi(h, a) + \phi \mathbb{E}_{h', a' \sim \pi(\cdot | h')} [Q^\pi(h', a')]\end{aligned}$$

where  $h'$  is the next history.

- Existing methods for estimating feature expectations fall into two extremes:
  - (1) model-based (online) IRL approaches learn a model of the world and then use the model as a simulator to obtain on-policy roll-outs
  - (2) batch IRL approaches use Q-learning for off-policy evaluation (Lee et al., 2019), and can only be used to evaluate policies similar to the expert policy and require warm start.

# Counterfactual $\mu$ - Learning

- The paper proposes counterfactual  $\mu$ -learning, a novel method for estimating feature expectations that uses these counterfactuals as part of temporal difference learning with 1-step bootstrapping.
- This approach falls in-between (1) and (2) and allows us to estimate feature expectations for any candidate policy  $\pi$  in the batch IRL setting.
- The counterfactual  $\mu$ -learning algorithm learns the  $\mu$ -values for policy  $\pi$  iteratively by updating the current estimates of the  $\mu$ -values with the feature map plus the  $\mu$ -values obtained by following policy  $\pi$  in the new counterfactual history  $h' = (h, a, \mathbb{E}[Y[a]|h])$

$$\hat{\mu}^\pi \leftarrow \hat{\mu}^\pi(h, a) + \alpha(\phi(h, a) + \gamma \mathbb{E}_{a' \sim \pi(\cdot|h')}[\hat{\mu}^\pi(h', a')] - \hat{\mu}^\pi(h', a'))$$

where  $\alpha$  is the learning rate.

# Batch, Max-Margin CIRL

---

**Algorithm 1** (Batch, Max-Margin) CIRL

---

- 1: **Input:** Batch dataset  $\mathcal{D}$ , max iterations  $n$ , convergence threshold  $\epsilon$ ,  
feature map  $\phi(h_t, a_t) = \mathbb{E}[Y_{t+1}[a_t]|h_t]$
  - 2:  $\mu^{\pi_E} \leftarrow$  compute  $\pi_E$ 's feature expectations (Equation 3)
  - 3:  $w_0 \leftarrow$  random initial reward weights,  $\pi_0 \leftarrow$  compute optimal policy for  $R_0 = w_0 \cdot \phi$
  - 4:  $\mu^{\pi_0} \leftarrow$  compute  $\pi_0$ 's feature expectations (counterfactual  $\mu$ -learning)
  - 5:  $\Pi = \{\pi_0\}, \Delta = \{\mu^{\pi_0}\}, \bar{\mu}_0 = \mu^{\pi_0}$
  - 6: **for**  $k = 1$  to  $n$  **do**
  - 7:    $w_k = \mu^{\pi_E} - \bar{\mu}_{k-1}$ ,  $\pi_k \leftarrow$  compute optimal policy for  $R_k = w_k \cdot \phi$
  - 8:    $\mu^{\pi_k} \leftarrow$  compute  $\pi_k$ 's feature expectations (counterfactual  $\mu$ -learning)
  - 9:    $\Pi = \Pi \cup \{\pi_k\}, \Delta = \Delta \cup \{\mu^{\pi_k}\}$
  - 10:   Orthogonally project  $\mu^{\pi_E}$  onto line through  $\bar{\mu}_{k-1}, \mu^{\pi_k}$ :  
      
$$\bar{\mu}_k = \frac{(\mu^{\pi_k} - \bar{\mu}_{k-1})^T (\mu^{\pi_E} - \bar{\mu}_{k-1})}{(\mu^{\pi_k} - \bar{\mu}_{k-1})^T (\mu^{\pi_k} - \bar{\mu}_{k-1})} (\mu^{\pi_k} - \bar{\mu}_{k-1}) + \bar{\mu}_{k-1}, \quad t = \|\mu^{\pi_E} - \bar{\mu}_k\|_2$$
  - 11:   **if**  $t < \epsilon$  **then break**
  - 12: **end for**
  - 13:  $K = \arg \min_{k; \mu^{\pi_k} \in \Delta} \|\mu^{\pi_E} - \mu^{\pi_k}\|_2, \tilde{R}(h, a) = w_K \cdot \phi(h, a)$
  - 14: **Output:**  $\tilde{R}, \Delta, \Pi$
- 

Figure: Psedo Code for CIRL

*Thank You*