



ALMA MATER STUDIORUM
UNIVERSITÀ DI BOLOGNA



POLITECNICO
MILANO 1863

Handling Non-Stationary Experts in Inverse Reinforcement Learning: A Water System Control Case Study

A. Likmeta

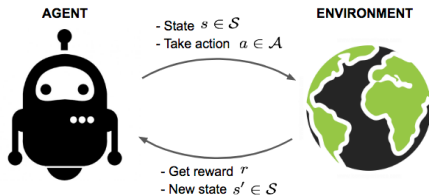
A. M. Metelli
M. Giuliani

G. Ramponi
M. Restelli

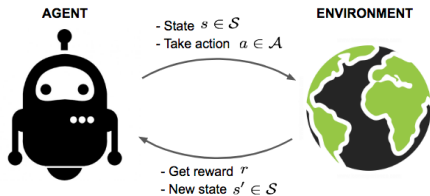
A. Tirinzoni

January 15, 2021
M2L Summer School

■ RL: Maximize sum of rewards

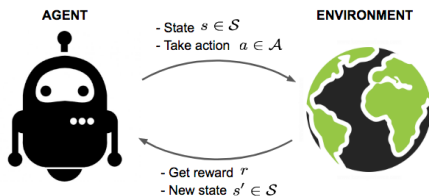


- **RL:** Maximize sum of rewards



! Reward function hard to define.

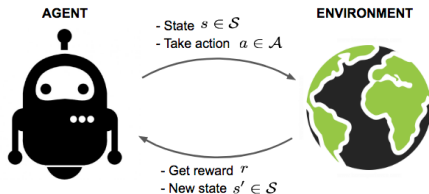
■ RL: Maximize sum of rewards



! Reward function hard to define.

! Demonstrations often available.

- **RL**: Maximize sum of rewards

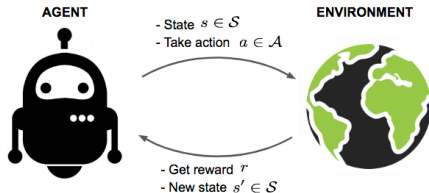


- ! Reward function hard to define.
- ! Demonstrations often available.

- **IRL**: Recovers **unknown reward** from expert demonstrations



- **RL**: Maximize sum of rewards



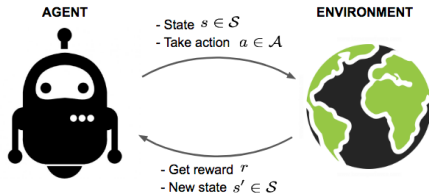
- ! Reward function hard to define.
- ! Demonstrations often available.

- **IRL**: Recovers **unknown reward** from expert demonstrations



- ! Using only demonstrations.

- **RL**: Maximize sum of rewards



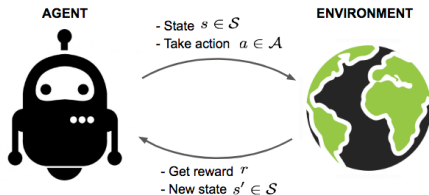
- ! Reward function hard to define.
- ! Demonstrations often available.

- **IRL**: Recovers **unknown reward** from expert demonstrations



- ! Using only demonstrations.
- ! No interactions with environment.

- **RL**: Maximize sum of rewards



- ! Reward function hard to define.
- ! Demonstrations often available.

- **IRL**: Recovers **unknown reward** from expert demonstrations



- ! Using only expert demonstrations.
- ! No interactions with environment.
- ! Non-stationarity in the objectives.

- MDP without Reward (MDP\(R) $\mathcal{M} = (\mathcal{S}, \mathcal{A}, P, \gamma, \mu)$ (Puterman, 1994)

- MDP without Reward (MDP\(R) $\mathcal{M} = (\mathcal{S}, \mathcal{A}, P, \gamma, \mu)$ (Puterman, 1994)
- Linearly parameterized rewards: $R_{\omega}(s, a) = \phi(s, a)^T \omega$

- MDP without Reward (MDP\(R) $\mathcal{M} = (\mathcal{S}, \mathcal{A}, P, \gamma, \mu)$ (Puterman, 1994)
- Linearly parameterized rewards: $R_{\omega}(s, a) = \phi(s, a)^T \omega$
- Parametrized policy π_{θ} :

$$J_{\pi}(\theta, \omega) = \mathbb{E}_{\substack{s_0 \sim \mu \\ a_t \sim \pi_{\theta}}} \left[\sum_{t=0}^{+\infty} \gamma^t R_{\omega}(s_t, a_t) \right] = \omega^T \psi(\theta)$$

- 1 Infer reward functions only from demonstrations
- 2 Handle Non-Stationarity in demonstrations
- 3 Use Case: Como Lake Water Control System

- 1 Infer reward functions only from demonstrations
- 2 Handle Non-Stationarity in demonstrations
- 3 Use Case: Como Lake Water Control System

- No interaction with the environment and no forward learning

- No interaction with the environment and no forward learning

Idea: if the expert's policy is optimal then $\nabla_{\theta} J(\theta, \omega) = 0$

- No interaction with the environment and no forward learning

Idea: if the expert's policy is optimal then $\nabla_{\theta} J(\theta, \omega) = 0$

Σ -Gradient Inverse Reinforcement Learning (Ramponi et al., 2020)

$$\min_{\|\omega\|_1=1} \left\| \hat{\nabla}_{\theta} \psi(\theta) \omega \right\|^2 \left[(\omega \otimes \mathbf{I}_d)^T \Sigma (\omega \otimes \mathbf{I}_d) \right]^{-1}$$

- No interaction with the environment and no forward learning

Idea: if the expert's policy is optimal then $\nabla_{\theta} J(\theta, \omega) = 0$

Σ -Gradient Inverse Reinforcement Learning (Ramponi et al., 2020)

$$\min_{\|\omega\|_1=1} \left\| \hat{\nabla}_{\theta} \psi(\theta) \omega \right\|^2_{\left[(\omega \otimes \mathbf{I}_d)^T \Sigma (\omega \otimes \mathbf{I}_d) \right]^{-1}}$$

- ! Need the parametrized expert policy (perform Behavioral Cloning).

- No interaction with the environment and no forward learning

Idea: if the expert's policy is optimal then $\nabla_{\theta} J(\theta, \omega) = 0$

Σ -Gradient Inverse Reinforcement Learning (Ramponi et al., 2020)

$$\min_{\|\omega\|_1=1} \left\| \hat{\nabla}_{\theta} \psi(\theta) \omega \right\|_{\left[(\omega \otimes \mathbf{I}_d)^T \Sigma (\omega \otimes \mathbf{I}_d) \right]^{-1}}^2$$

- ! Need the parametrized expert policy (perform Behavioral Cloning).
- ! Gradients are estimated from demonstrations
→ could not exist a ω s.t. $\hat{\nabla}_{\theta} J(\theta, \omega) = 0$

- No interaction with the environment and no forward learning

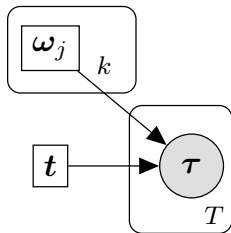
Idea: if the expert's policy is optimal then $\nabla_{\theta} J(\theta, \omega) = 0$

Σ -Gradient Inverse Reinforcement Learning (Ramponi et al., 2020)

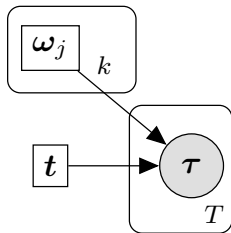
$$\min_{\|\omega\|_1=1} \left\| \hat{\nabla}_{\theta} \psi(\theta) \omega \right\|^2_{\left[(\omega \otimes \mathbf{I}_d)^T \Sigma (\omega \otimes \mathbf{I}_d) \right]^{-1}}$$

- ! Need the parametrized expert policy (perform Behavioral Cloning).
- ! Gradients are estimated from demonstrations
 - could not exist a ω s.t. $\hat{\nabla}_{\theta} J(\theta, \omega) = 0$
 - account for the uncertainty in the Jacobian estimation
 - allows computing likelihood of D given ω , $p(D|\omega)$

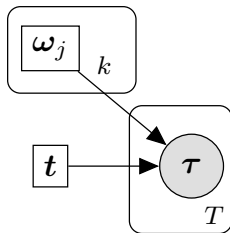
- 1 Infer reward functions only from demonstrations
- 2 Handle Non-Stationarity in demonstrations
- 3 Use Case: Como Lake Water Control System



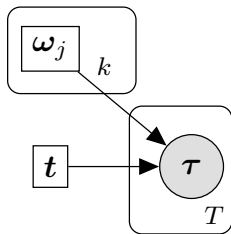
- Consider dataset $D = \{\tau_i\}_{i=1}^T$ as a lifelong trajectory $\tau = (\tau_1 | \dots | \tau_T)$



- Consider dataset $D = \{\tau_i\}_{i=1}^T$ as a lifelong trajectory $\tau = (\tau_1 | \dots | \tau_T)$
- Non-stationary behavior characterized by $k \leq T$ regimes



- Consider dataset $D = \{\tau_i\}_{i=1}^T$ as a lifelong trajectory $\tau = (\tau_1 | \dots | \tau_T)$
- Non-stationary behavior characterized by $k \leq T$ regimes
- Find change points
 $1 < t_1 < \dots < t_{k-1} < t_k = T$



- Consider dataset $D = \{\tau_i\}_{i=1}^T$ as a lifelong trajectory $\tau = (\tau_1 | \dots | \tau_T)$
- Non-stationary behavior characterized by $k \leq T$ regimes
- Find change points
 $1 < t_1 < \dots < t_{k-1} < t_k = T$
- Find rewards $\mathbf{R} = (R_{\omega_1}, \dots, R_{\omega_k})$

Idea: Take inspiration from *change-point detection* (Aminikhanghahi and Cook, 2017)

Idea: Take inspiration from *change-point detection* (Aminikhanghahi and Cook, 2017)

NS- Σ -GIRL (adaptation of *Opt* (Truong et al., 2020))

- 1 Define Likelihood of τ under solution $\Omega = (\omega_1, \dots, \omega_k, t_1, \dots, t_{k-1})$

$$\mathcal{L}(\Omega|\tau) = p(\tau|\Omega) = \prod_{i=1}^T \sum_{j=1}^k p(\tau_i|\omega_j) \mathbb{1}_{\{i \in I_j\}}$$

- 2 For each $1 \leq u < v \leq T$ solve:

$$\min_{\substack{\omega_{uv} \in \mathbb{R}_{\geq 0}^q \\ \|\omega_{uv}\|_1 = 1}} (v-u) \sum_{i=u}^{v-1} \left\| \hat{\nabla}_{\theta} \psi_i(\theta) \omega_{uv} \right\|_{[(\omega_{uv} \otimes \mathbf{I}_d) \Sigma_i (\omega_{uv} \otimes \mathbf{I}_d)^T]^{-1}}^2.$$

- 1 Infer reward functions only from demonstrations
- 2 Handle Non-Stationarity in demonstrations
- 3 Use Case: Como Lake Water Control System

- **Problem:** Infer operator intentions from hystorical dam operation.

- **Problem:** Infer operator intentions from hystorical dam operation.
- **Problem:** The intentions of the operators change during 60 years

- **Problem:** Infer operator intentions from hystorical dam operation.
- **Problem:** The intentions of the operators change during 60 years
- Continuous state: water stored in the lake S_t , a continuous action: water released a_t , a state-transition function of: lake inflow q_{t+1}

$$S_{t+1} = S_t + q_{t+1} - r_{t+1}(S_t, a_t, q_{t+1})$$

- **Problem:** Infer operator intentions from hystorical dam operation.
- **Problem:** The intentions of the operators change during 60 years
- Continuous state: water stored in the lake S_t , a continuous action: water released a_t , a state-transition function of: lake inflow q_{t+1}

$$S_{t+1} = S_t + q_{t+1} - r_{t+1}(S_t, a_t, q_{t+1})$$

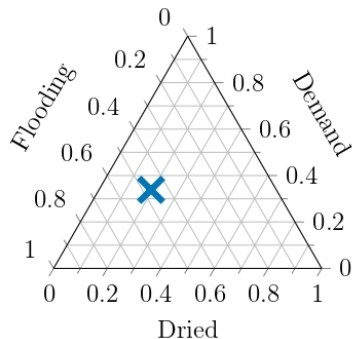
- Data are daily and were provided by Consorzio dell'Adda (www.addaconsorzio.it)

- Three reward features representing conflicting objectives:

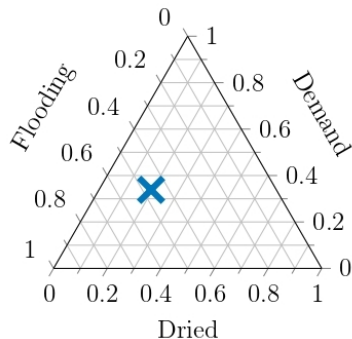
- Three reward features representing conflicting objectives:
- *Supply deficit* - (ϕ^D): deficit between the release and the demand

- Three reward features representing conflicting objectives:
- *Supply deficit* - (ϕ^D): deficit between the release and the demand
- *Flood risk* (ϕ^F): penalize small releases associated to high lake levels

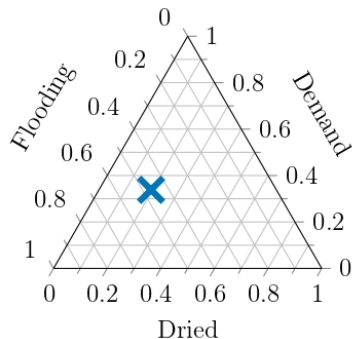
- Three reward features representing conflicting objectives:
- *Supply deficit* - (ϕ^D): deficit between the release and the demand
- *Flood risk* (ϕ^F): penalize small releases associated to high lake levels
- *Drought risk* (ϕ^L): penalize large releases with low lake levels



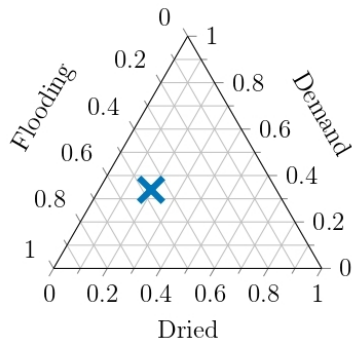
- Slight predominance for controlling the floods ($\omega^F = 0.47$)



- Slight predominance for controlling the floods ($\omega^F = 0.47$)
- Remaining weight is divided between demand ($\omega^D = 0.34$) and drought control ($\omega^L = 0.19$)



- Slight predominance for controlling the floods ($\omega^F = 0.47$)
- Remaining weight is divided between demand ($\omega^D = 0.34$) and drought control ($\omega^L = 0.19$)
- Results in line with literature results (Giuliani et al., 2019)



- Slight predominance for controlling the floods ($\omega^F = 0.47$)
- Remaining weight is divided between demand ($\omega^D = 0.34$) and drought control ($\omega^L = 0.19$)
- Results in line with literature results (Giuliani et al., 2019)
- Expert almost Pareto optimal



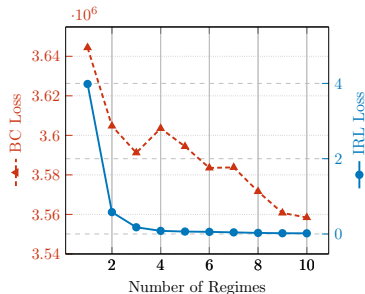
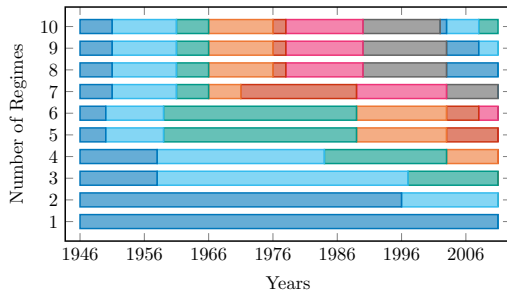
Recent Flooding (Como, 2020)

- Slight predominance for controlling the floods ($\omega^F = 0.47$)
- Remaining weight is divided between demand ($\omega^D = 0.34$) and drought control ($\omega^L = 0.19$)
- Results in line with literature results (Giuliani et al., 2019)
- Expert almost Pareto optimal

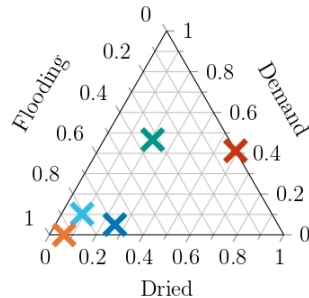
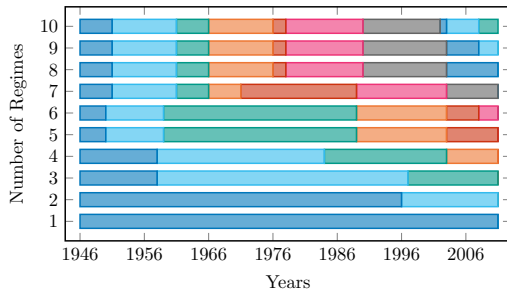
- **Problem:** Lake Como is a **non-stationary** system that has undergone several alterations

- **Problem:** Lake Como is a **non-stationary** system that has undergone several alterations
- **Idea:** Consider the dataset as a lifelong trajectory. Make subdivisions yearly. Find K **intention change points**.

- **Problem:** Lake Como is a **non-stationary** system that has undergone several alterations
- **Idea:** Consider the dataset as a lifelong trajectory. Make subdivisions yearly. Find K **intention change points**.



- **Problem:** Lake Como is a **non-stationary** system that has undergone several alterations
- **Idea:** Consider the dataset as a lifelong trajectory. Make subdivisions yearly. Find K **intention change points**.



- Extended Σ -GIRL to the non-stationary setting
- Carefully designed reward features that explain the expert behaviour
- Identified several hystorical regimes of interest

Thank You for Your Attention!

- Samaneh Aminikhanghahi and Diane J Cook. A survey of methods for time series change point detection. *Knowledge and information systems*, 51(2):339–367, 2017.
- Qui Como. Il lago è fuori: le nuove barriere anti-esondazione, October 2020. URL <https://www.quicomo.it/attualita/lago-di-como-fuori-3-ottobre-2020.html>.
- Matteo Giuliani, Marta Zaniolo, Andrea Castelletti, Guido Davoli, and Paul Block. Detecting the state of the climate system via artificial intelligence to improve seasonal forecasts and inform reservoir operations. *Water Resources Research*, 55:9133–9147, 2019.
- Martin L. Puterman. *Markov Decision Processes: Discrete Stochastic Dynamic Programming*. John Wiley & Sons, Inc., New York, NY, USA, 1994. ISBN 0471619779.
- Giorgia Ramponi, Amarildo Likmeta, Alberto Maria Metelli, Andrea Tirinzoni, and Marcello Restelli. Truly batch model-free inverse reinforcement learning about multiple intentions. In Silvia Chiappa and Roberto Calandra, editors, *Proceedings of the Twenty Third International Conference on Artificial Intelligence and Statistics*, volume 108 of *Proceedings of Machine Learning Research*, pages 2359–2369, Online, 26–28 Aug 2020. PMLR. URL <http://proceedings.mlr.press/v108/ramponi20a.html>.
- Charles Truong, Laurent Oudre, and Nicolas Vayatis. Selective review of offline change point detection methods. *Signal Processing*, 167:107299, 2020.