

PROBLEM AND MOTIVATION

- **Reinforcement Learning** (RL, Sutton and Barto, 2018): find optimal policy π^* maximizing the *expected return* $J^\pi(\theta, \omega)$:

$$J_\pi(\theta, \omega) = \mathbb{E}_{\substack{S_0 \sim \mu \\ A_t \sim \pi_\theta(\cdot | S_t) \\ S_{t+1} \sim \mathcal{P}(\cdot | S_t, A_t)}} \left[\sum_{t=0}^{+\infty} \gamma^t R_\omega(S_t, A_t) | S_0 = s \right]$$

- **Inverse Reinforcement Learning**

1. Given a dataset D of demonstrations from an expert, find the unknown reward function being optimized:

$$R_{\pi^E}^* \in \left\{ R \in \mathcal{R} : \pi^E \in \arg \max_{\pi \in \Pi} J(\pi, R) \right\}.$$

2. We further assume linearity of the reward function in terms of a feature function ϕ

$$\mathcal{R} = \left\{ R_\omega = \omega^T \phi : \omega \in \mathbb{R}_{\geq 0}^q, \|\omega\|_1 = 1 \right\}$$

Model-Free No model of the environment is available

Batch We cannot further interact with the environment

GRADIENT INVERSE REINFORCEMENT LEARNING

- if π_{θ^E} is optimal for the reward R_{ω^E} , θ^E is a **stationary point** of the return $J(\theta, \omega^E) = (\omega^E)^T \psi(\theta)$
 $\implies \nabla_{\theta} J(\theta^E, \omega^E) = \nabla_{\theta} \psi(\theta^E) \omega^E = 0$
- Find the reward weights ω that lie in the **null space** of the jacobian
 \implies Due to **estimation error** the null space of the jacobian might be **empty**

CONTRIBUTIONS

- **Jacobian correction**: Account for the **uncertainty**, by modelling the sample distribution $\hat{\nabla}_{\theta} \psi(\theta) \sim \mathcal{N}(\mathbf{M}, \frac{1}{n} \Sigma)$

$$\min_{\substack{\omega \in \mathbb{R}_{\geq 0}^q \\ \|\omega\|_1 = 1}} \left\| \hat{\nabla}_{\theta} \psi(\theta) \omega \right\|_{[(\omega \otimes \mathbf{I}_d)^T \Sigma (\omega \otimes \mathbf{I}_d)]^{-1}}^2, (\Sigma\text{-GIRL})$$

- **Non-Stationarity**: Account for **non-stationary** behavior in the demonstrations due to **changing intentions**
 - **Detect** K intention change points and **identify** the reward functions $\{\omega\}_{j=1}^K$

$$\min_{\substack{\omega_{uv} \in \mathbb{R}_{\geq 0}^q \\ \|\omega_{uv}\|_1 = 1}} (v - u) \sum_{i=u}^{v-1} \left\| \hat{\nabla}_{\theta} \psi_i(\theta) \omega_{uv} \right\|_{[(\Sigma(u,v))_i]^{-1}}^2$$

- **Real-World Application**: Application of IRL in real-world dataset of Lake Como dam operation

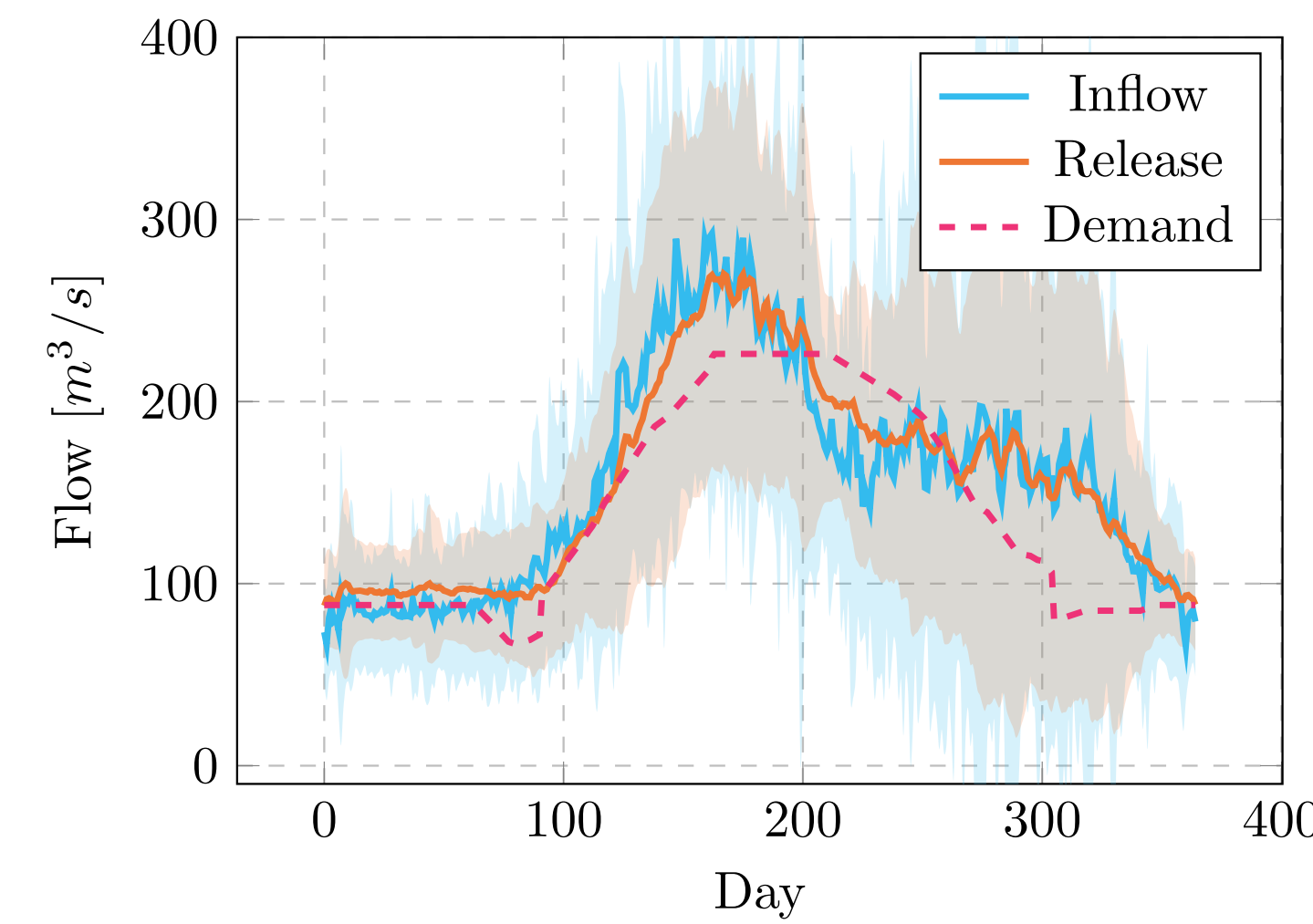
COMO USE CASE

System Modelling

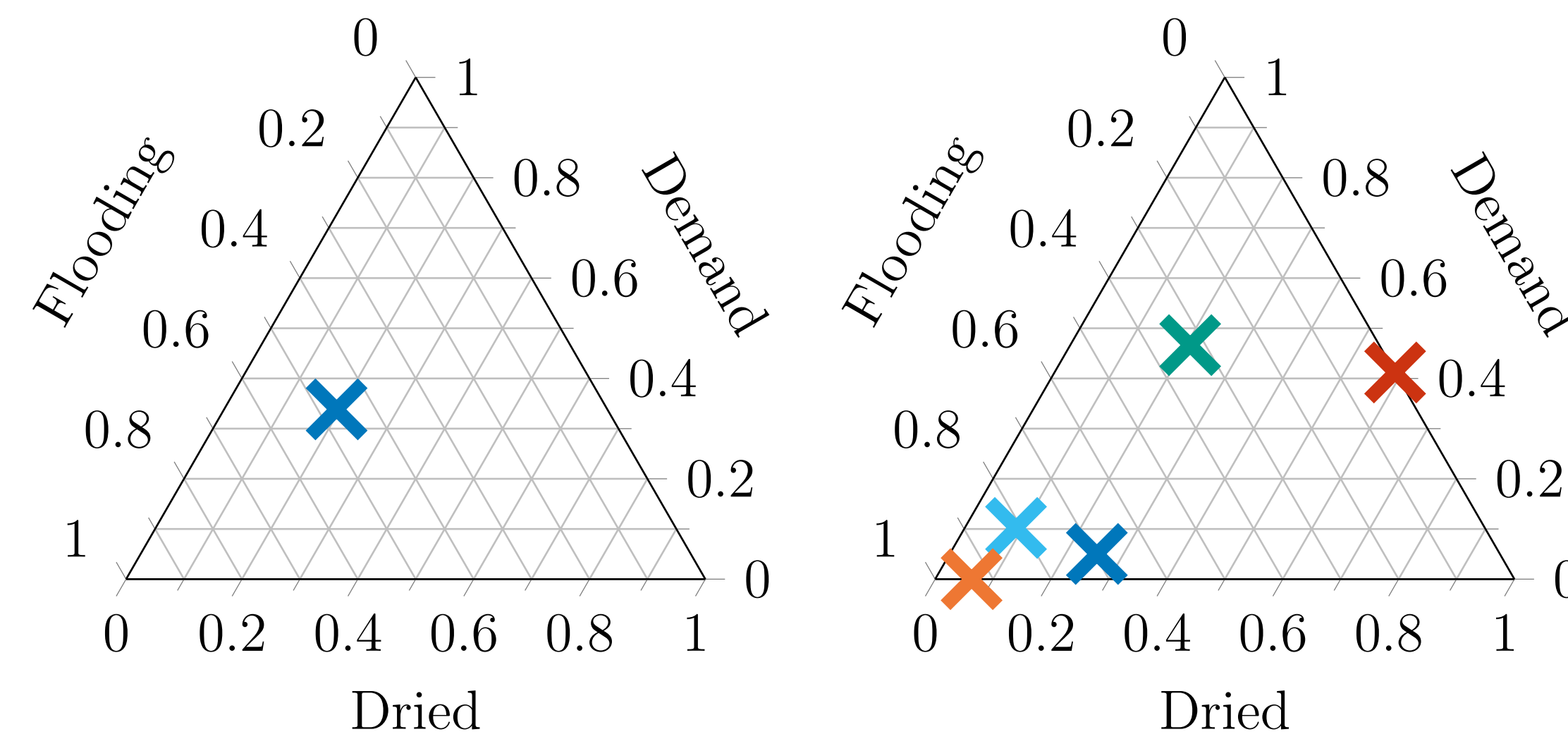
- **Problem**: Infer operator intentions from hystorical dam operation.
- **Problem**: The intentions of the operators change during 60 years
- Model as discrete-time, periodic, nonlinear, stochastic MDP
- Continuous state: water stored in the lake S_t , a continuous action: water released a_t , a state-transition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

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

- 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



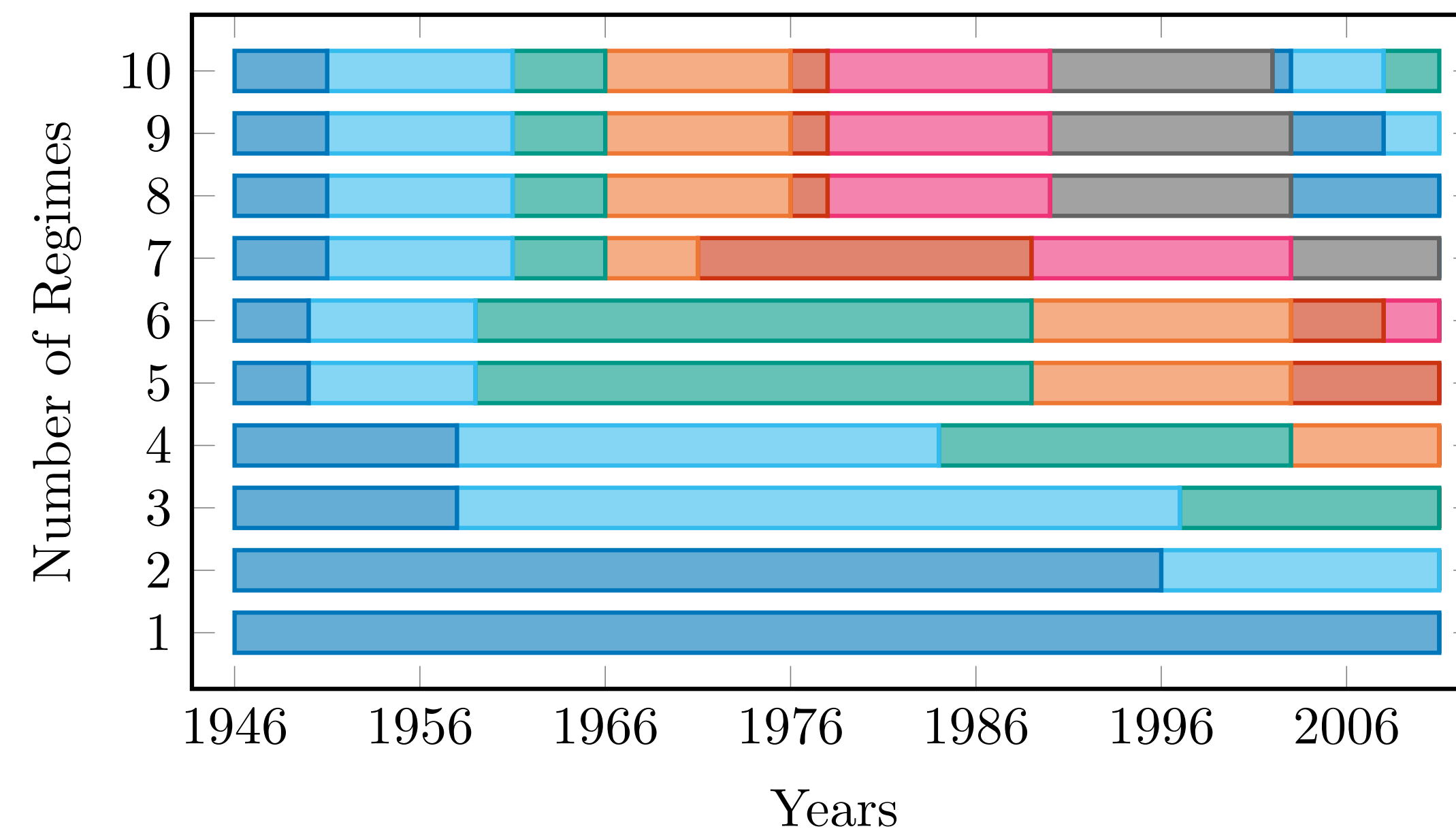
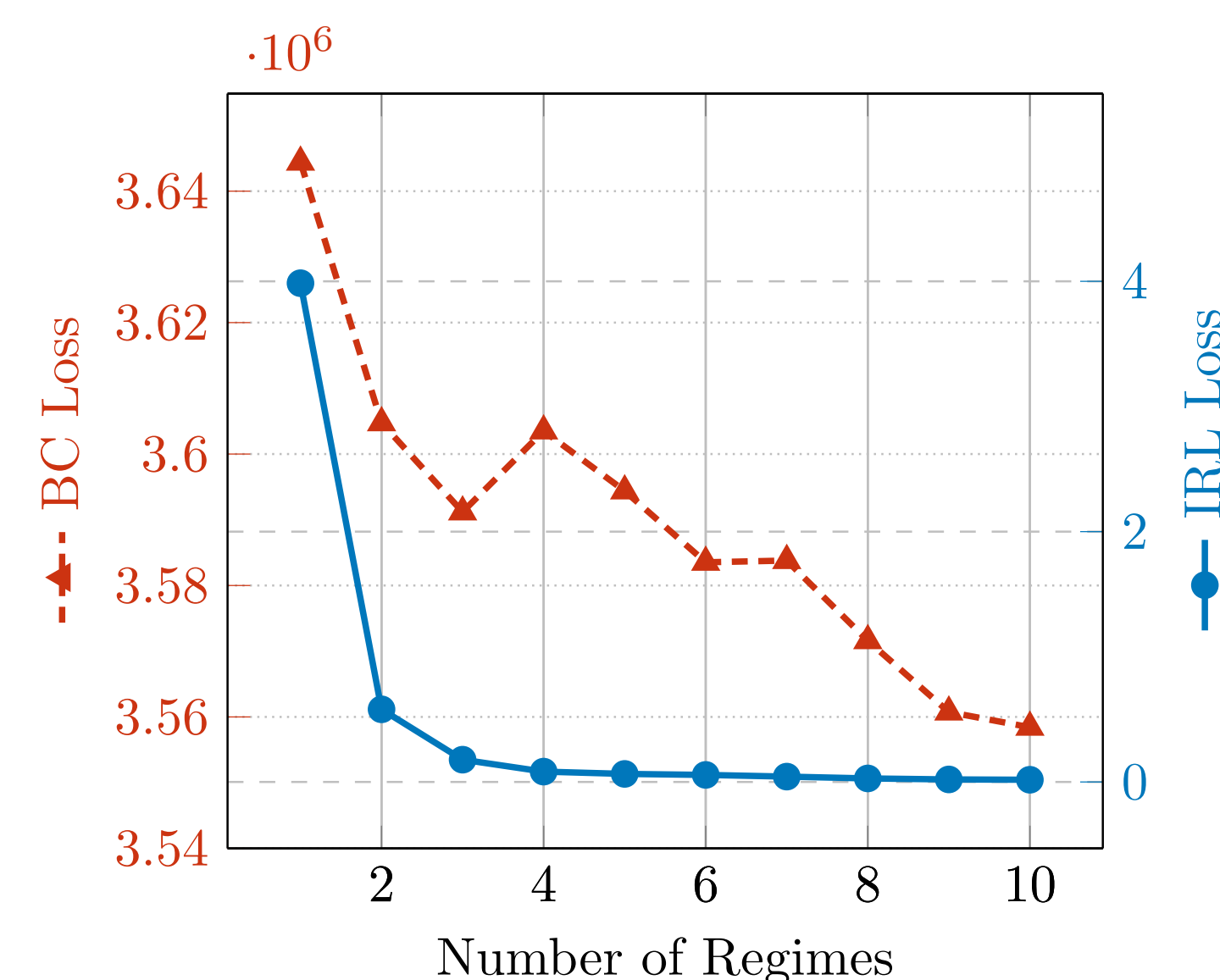
Single Reward IRL Results



- Slight predominance of the interest in 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

Non-Stationary IRL Results

- **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**.



- From an elbow analysis we can identify 4 or 5 distinct periods.
- Intervals match historical events such as, big flooding events (1987), or extreme droughts (2003)

Σ -GIRL (RAMPONI ET AL., 2020)

input: dataset of demonstrations D

output: optimal parameters ω^*

1. Perform BC to find the policy parameters θ^E
2. Estimate $\hat{\nabla}_{\theta} \psi(\theta)$ and Σ from the demonstration dataset
3. Find the weights ω^* minimizing

$$\omega^* = \arg \min \left\| \hat{\nabla}_{\theta} \psi(\theta) \omega \right\|_{[(\omega \otimes \mathbf{I}_d)^T \Sigma (\omega \otimes \mathbf{I}_d)]^{-1}}^2$$
4. **return** ω^*

NS- Σ -GIRL (LIKMETA ET AL., 2020)

input: dataset $\tau = (\tau_1 | \tau_2 | \dots | \tau_T)$, number of regimes k
output: optimal parameters $\Omega = (\omega_1, \dots, \omega_k, t_1, \dots, t_{k-1})$

1. **for** $u = 1, \dots, T-1$ **do**
2. **for** $v = u+1, \dots, T$ **do**
3. Define $D_{uv} = \{\tau_u, \dots, \tau_v-1\}$
4. Perform BC to find the policy parameters θ_{uv}
5. Optimize $\omega_{uv}^* \in \arg \max_{\omega} \log p(D_{uv} | \omega)$
6. $C_1(u, v) = \log p(D_{uv} | \omega_{uv}^*)$
7. **end for**
8. **end for**
9. **for** $l = 2, \dots, k-1$ **do**
10. **for** $u = 1, \dots, T-l$ **do**
11. **for** $v = u+l, \dots, T$ **do**
12. $C_l(u, v) = \max_{u+l-1 \leq t < v} \{C_{l-1}(u, t) + C_1(t+1, v)\}$
13. **end for**
14. **end for**
15. **end for**
16. $t_k = T$
17. **for** $l = k, \dots, 2$ **do**
18. $t_{l-1} \in \arg \max_{l-1 \leq t < t_l} C_{l-1}(1, t) + C_1(t+1, t_l)$
19. $\omega_l = \omega_{t_{l-1} t_l}^*$
20. **end for**
21. **return** Ω

REFERENCES

- M. Giuliani, M. Zaniolo, A. Castelletti, G. Davoli, and P. 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.
- A. Likmeta, G. Ramponi, A. M. Metelli, M. Tirinzoni, Andrea Giuliani, and M. Restelli. Dealing with multiple experts and non-stationarity in inverse reinforcement learning: an application to real-life problems. In *Machine Learning Journal (Under Revision)*, 2020.
- G. Ramponi, A. Likmeta, A. M. Metelli, A. Tirinzoni, and M. Restelli. Truly batch model-free inverse reinforcement learning about multiple intentions. In *The 23rd International Conference on Artificial Intelligence and Statistics*, 2020.
- R. S. Sutton and A. G. Barto. *Reinforcement learning: An introduction*. MIT press, 2018.