

The Importance of Pessimism in Fixed-Dataset Policy Optimization

Jacob Buckman (Mila; McGill University)*

Carles Gelada (OpenAI)

Marc G. Bellemare (Google Research; Mila; McGill University; CIFAR Fellow)

Our goal is an algo with minimal worst-case suboptimality,

$$\text{SUBOPT}(\mathcal{O}(D)) = \mathbb{E}_\rho[\mathbf{v}_{\mathcal{M}}^{\pi^*}] - \mathbb{E}_\rho[\mathbf{v}_{\mathcal{M}}^{\mathcal{O}(D)}].$$

We consider “value-based” algorithms,

$$\mathcal{O}_{\text{sub}}^{\text{VB}}(D) := \arg \max_{\pi} \mathbb{E}_\rho[\mathcal{E}_{\text{sub}}(D, \pi)].$$

These algorithms can be characterized by the choice of fixed point of \mathbb{E}_{sub} . Suboptimality of these algos permits an “over/under decomposition”,

Theorem 1. For any space \mathcal{X} , objective $f : \mathcal{X} \rightarrow \mathbb{R}$, and proxy objective $\hat{f} : \mathcal{X} \rightarrow \mathbb{R}$,

$$f(x^*) - f(\hat{x}^*) \leq \inf_{x \in \mathcal{X}} \left([f(x^*) - f(x)] + [f(x) - \hat{f}(x)] \right) + \sup_{x \in \mathcal{X}} \left(\hat{f}(x) - f(x) \right)$$

where $x^* := \arg \max_{x \in \mathcal{X}} f(x)$ and $\hat{x}^* := \arg \max_{x \in \mathcal{X}} \hat{f}(x)$. Furthermore, this bound is tight.

One important type of algo is “naive”:

$$f_{\text{naive}}(\mathbf{v}^\pi) := A^\pi(\mathbf{r}_D + \gamma P_D \mathbf{v}^\pi). \quad \text{SUBOPT}(\mathcal{O}_{\text{naive}}^{\text{VB}}(D)) \leq \inf_{\pi} \left(\mathbb{E}_\rho[\mathbf{v}_{\mathcal{M}}^{\pi^*} - \mathbf{v}_{\mathcal{M}}^\pi] + \mathbb{E}_\rho[\boldsymbol{\mu}_{D,\delta}^\pi] \right) + \sup_{\pi} \mathbb{E}_\rho[\boldsymbol{\mu}_{D,\delta}^\pi]$$

This often leads to a large “sup” term. We can fix this by finding pessimistic fixed points, which let us choose the relative size of the two terms:

$$\text{SUBOPT}(\mathcal{O}_{\text{ua}}^{\text{VB}}(D)) \leq \inf \left(\mathbb{E}_\rho[\mathbf{v}_{\mathcal{M}}^{\pi^*} - \mathbf{v}_{\mathcal{M}}^\pi] + (1 + \alpha) \cdot \mathbb{E}_\rho[\boldsymbol{\mu}_{D,\delta}^\pi] \right) + (1 - \alpha) \cdot \left(\sup_{\pi} \mathbb{E}_\rho[\boldsymbol{\mu}_{D,\delta}^\pi] \right)$$

$$f_{\text{ua}}(\mathbf{v}^\pi) = A^\pi(\mathbf{r}_D + \gamma P_D \mathbf{v}^\pi) - \alpha \mathbf{u}_{D,\delta}^\pi$$

Implementing this algorithm requires implementing a valid uncertainty measure, which we don’t know how to do right now with NNs. If we take “trivial uncertainty” of V_{max} , we get proximal algorithms:

$$f_{\text{proximal}}(\mathbf{v}^\pi) = A^\pi(\mathbf{r}_D + \gamma P_D \mathbf{v}^\pi) - \alpha \left(\frac{TV_S(\pi, \hat{\pi}_D)}{(1 - \gamma)^2} \right)$$

The trivial uncertainty is the “worst” uncertainty, leading to a much looser bound; but it is, at least, implementable.

This work **provides formal justification** for the properties of **every “Offline RL” algorithm** in the literature, including:

BCQ, CRR, SPIBB, BEAR, CQL, KLC, BRAC, MBS-QI, MoREL, MOPO, and more.

