

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”,

$$\text{SUBOPT}(\mathcal{O}^{\text{VB}}(D)) \leq \inf_{\pi} \left(\mathbb{E}_\rho[\mathbf{v}_{\mathcal{M}}^{\pi^*} - \mathbf{v}_{\mathcal{M}}^{\pi}] + \mathbb{E}_\rho[\mathbf{v}_{\mathcal{M}}^{\pi} - \mathbf{v}_D^{\pi}] \right) + \sup_{\pi} \left(\mathbb{E}_\rho[\mathbf{v}_D^{\pi} - \mathbf{v}_{\mathcal{M}}^{\pi}] \right)$$

The actual suboptimality depends choice of \mathbb{E}_{sub} . 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:

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

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

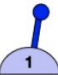


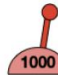
Implementing this algorithm requires implementing a valid uncertainty measure, which we don’t know how to do right now with NNs. If we take a “trivial uncertainty” of V_{max} , we get proximal algorithms:

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

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.

				...	
	1	2	3		1000
μ	99%	1%	1%		1%
n_D	10000	1	1		1
μ_D	98%	0%	100%		0%

