
Supplementary Material for Probability Functional Descent

A. Proofs and Computations

Lemma 1. *Let $J : \mathcal{P}(X) \rightarrow \mathbb{R}$. Then $\Psi : X \rightarrow \mathbb{R}$ is an influence function of J at μ if and only if*

$$\left. \frac{d}{d\epsilon} J(\mu + \epsilon\chi) \right|_{\epsilon=0^+} = \int_X \Psi(x) \chi(dx).$$

Proof. The left-hand side equals (1), which equals (2). \square

Theorem 1 (Chain rule). *Let $J : \mathcal{P}(X) \rightarrow \mathbb{R}$ be continuously differentiable, in the sense that the influence function Ψ_μ exists and $(\mu, \nu) \mapsto \mathbb{E}_{x \sim \nu}[\Psi_\mu(x)]$ is continuous. Let the parameterization $\theta \mapsto \mu_\theta$ be differentiable, in the sense that $\frac{1}{\|h\|}(\mu_{\theta+h} - \mu_\theta)$ converges to a weak limit as $h \rightarrow 0$. Then*

$$\nabla_\theta J(\mu_\theta) = \nabla_\theta \mathbb{E}_{x \sim \mu_\theta}[\hat{\Psi}(x)],$$

where $\hat{\Psi} = \Psi_{\mu_\theta}$ is treated as a function $X \rightarrow \mathbb{R}$ that is not dependent on θ .

Proof. Without loss of generality, assume $\theta \in \mathbb{R}$, as the gradient is simply a vector of one-dimensional derivatives. Let $\chi_\epsilon = \frac{1}{\epsilon}(\mu_{\theta+\epsilon} - \mu_\theta)$, and let $\chi = \lim_{\epsilon \rightarrow 0} \chi_\epsilon$ (weakly). Then

$$\begin{aligned} \frac{d}{d\theta} J(\mu_\theta) &= \left. \frac{d}{d\epsilon} J(\mu_{\theta+\epsilon}) \right|_{\epsilon=0} \\ &= \left. \frac{d}{d\epsilon} J(\mu_\theta + \epsilon\chi_\epsilon) \right|_{\epsilon=0}. \end{aligned}$$

Assuming for now that

$$\left. \frac{d}{d\epsilon} J(\mu_\theta + \epsilon\chi_\epsilon) \right|_{\epsilon=0} = \left. \frac{d}{d\epsilon} J(\mu_\theta + \epsilon\chi) \right|_{\epsilon=0},$$

we have by Lemma 1 that

$$\begin{aligned} \frac{d}{d\theta} J(\mu_\theta) &= \int_X \hat{\Psi} d\chi \\ &= \int_X \hat{\Psi} d\left(\lim_{\epsilon \rightarrow 0} \frac{1}{\epsilon}(\mu_{\theta+\epsilon} - \mu_\theta)\right) \\ &= \lim_{\epsilon \rightarrow 0} \int_X \hat{\Psi} d\left(\frac{1}{\epsilon}(\mu_{\theta+\epsilon} - \mu_\theta)\right) \\ &= \frac{d}{d\theta} \int_X \hat{\Psi} d\mu_\theta, \end{aligned}$$

where the interchange of limits is by the definition of weak convergence (recall we assumed that X is compact, so $\hat{\Psi}$ is continuous and bounded by virtue of being continuous).

The equality we assumed is the definition of a stronger notion of differentiability called Hadamard differentiability of J . Our conditions imply Hadamard differentiability via Proposition 2.33 of Penot (2012), noting that the map $(\mu, \chi) \mapsto \int_X \Psi_\mu d\chi$ is continuous by assumption. \square

Theorem 2 (Fenchel–Moreau representation). *Let $J : \mathcal{M}(X) \rightarrow \overline{\mathbb{R}}$ be proper, convex, and lower semicontinuous. Then the maximizer of $\varphi \mapsto \mathbb{E}_{x \sim \mu}[\varphi(x)] - J^*(\varphi)$, if it exists, is an influence function for J at μ . With some abuse of notation, we have that*

$$\Psi_\mu = \arg \max_{\varphi \in \mathcal{C}(X)} \left[\mathbb{E}_{x \sim \mu}[\varphi(x)] - J^*(\varphi) \right].$$

Proof. We will exploit the Fenchel–Moreau theorem, which applies in the setting of locally convex, Hausdorff topological vector spaces (see e.g. Zalinescu (2002)). The space we consider is $\mathcal{M}(X)$, the space of signed, finite measures equipped with the topology of weak convergence, of which $\mathcal{P}(X)$ is a convex subset. $\mathcal{M}(X)$ is indeed locally convex and Hausdorff, and its dual space is $\mathcal{C}(X)$ (see e.g. Aliprantis & Border (2006), section 5.14).

We now show that a maximizer φ^* is an influence function. By the Fenchel–Moreau theorem,

$$J(\mu) = J^{**}(\mu) = \sup_{\varphi \in \mathcal{C}(X)} \left[\int_X \varphi d\mu - J^*(\varphi) \right],$$

and

$$J(\mu + \epsilon\chi) = \sup_{\varphi \in \mathcal{C}(X)} \left[\int_X \varphi d\mu + \epsilon \int_X \varphi d\chi - J^*(\varphi) \right].$$

Because J is differentiable, $\epsilon \mapsto J(\mu + \epsilon\chi)$ is differentiable, so by the envelope theorem (Milgrom & Segal, 2002),

$$\left. \frac{d}{d\epsilon} J(\mu + \epsilon\chi) \right|_{\epsilon=0} = \int_X \varphi^* d\chi,$$

so that φ^* is an influence function by Lemma 1.

The abuse of notation stems from the fact that not all influence functions are maximizers. This is true, though, if

$J(\mu) = \infty$ if $\mu \notin \mathcal{P}(X)$:

$$\begin{aligned} & \int_X \Psi_\mu d\mu - J^*(\Psi_\mu) \\ &= \int_X \Psi_\mu d\mu - \sup_{\nu \in \mathcal{P}(X)} \left[\int_X \Psi_\mu d\nu - J(\nu) \right] \\ &= \inf_{\nu \in \mathcal{P}(X)} \left[- \int_X \Psi_\mu d(\nu - \mu) + J(\nu) \right] \\ &= \inf_{\nu \in \mathcal{P}(X)} \left[- \frac{d}{d\epsilon} J(\mu + \epsilon(\nu - \mu)) \Big|_{\epsilon=0} + J(\nu) \right] \\ &\geq J(\mu), \end{aligned}$$

since the convex function $f(\epsilon) = J(\mu + \epsilon(\nu - \mu))$ lies above its tangent line:

$$f(1) \geq f(0) + 1 \cdot f'(0).$$

Since $J(\mu) = J^{**}(\mu)$, we have that

$$\int_X \Psi_\mu d\mu - J^*(\Psi_\mu) \geq \sup_{\varphi \in \mathcal{C}(X)} \left[\int_X \varphi d\mu - J^*(\varphi) \right].$$

□

The following lemma will come in handy in our computations.

Lemma 2. Suppose $J : \mathcal{M}(X) \rightarrow \overline{\mathbb{R}}$ has a representation

$$J(\mu) = \sup_{\varphi \in \mathcal{C}(X)} \left[\int_X \varphi d\mu - K(\varphi) \right],$$

where $K : \mathcal{C}(X) \rightarrow \overline{\mathbb{R}}$ is proper, convex, and lower semi-continuous. Then $J^* = K$.

Proof. By definition of the convex conjugate, $J = K^*$. Then $J^* = K^{**} = K$, by the Fenchel–Moreau theorem. □

We note that when applying this lemma, we will often implicitly define the appropriate extension of J to $\mathcal{M}(X)$ to be $J(\mu) = \sup_{\varphi \in \mathcal{C}(X)} [\int \varphi d\mu - K(\varphi)]$. The exact choice of extension can certainly affect the exact form of the convex conjugate; see Ruderman et al. (2012) for one example of this phenomenon.

Proposition 2. Suppose μ has density $p(x)$ and ν has density $q(x)$. Then the influence function for J_{JS} is

$$\Psi_{\text{JS}}(x) = \frac{1}{2} \log \frac{p(x)}{p(x) + q(x)}.$$

Proof. The result follows from Lemma 1:

$$\begin{aligned} & \frac{d}{d\epsilon} J_{\text{JS}}(\mu + \epsilon\chi) \Big|_{\epsilon=0} \\ &= \frac{1}{2} \int_X \frac{d}{d\epsilon} \left[(p + \epsilon\chi) \log \frac{p + \epsilon\chi}{\frac{1}{2}(p + \epsilon\chi) + \frac{1}{2}q} \right. \\ &\quad \left. + q \log \frac{q}{\frac{1}{2}(p + \epsilon\chi) + \frac{1}{2}q} \right]_{\epsilon=0} dx \\ &= \frac{1}{2} \int_X \left[\log \frac{p}{\frac{1}{2}p + \frac{1}{2}q} + 1 - \frac{p}{p+q} - \frac{q}{p+q} \right] \chi dx \\ &= \frac{1}{2} \int_X \left[\log \frac{p}{p+q} + \log 2 \right] \chi dx. \end{aligned}$$

□

Proposition 3. The convex conjugate of J_{JS} is

$$J_{\text{JS}}^*(\varphi) = -\frac{1}{2} \mathbb{E}_{x \sim \nu} [\log(1 - e^{2\varphi(x) + \log 2})] - \frac{1}{2} \log 2.$$

Proof.

$$\begin{aligned} J_{\text{JS}}^*(\varphi) &= \sup_{\mu \in \mathcal{M}(X)} \left[\int_X \varphi d\mu - J_{\text{JS}}(\mu) \right] \\ &= \sup_p \int_X \left[\varphi p - \frac{1}{2} p \log \frac{p}{\frac{1}{2}p + \frac{1}{2}q} - \frac{1}{2} q \log \frac{q}{\frac{1}{2}p + \frac{1}{2}q} \right] dx. \end{aligned}$$

Setting the integrand's derivative w.r.t. p to 0, we find that pointwise, the optimal p satisfies

$$\varphi = \frac{1}{2} \log \frac{p}{\frac{1}{2}p + \frac{1}{2}q}.$$

We eliminate p in the integrand. Notice that the first two terms in the integrand cancel after plugging in p . Since

$$\frac{q}{\frac{1}{2}p + \frac{1}{2}q} = 2 \left(1 - \frac{p}{p+q} \right) = 2(1 - 2e^{2\varphi}),$$

we obtain that

$$J_{\text{JS}}^*(\varphi) = -\frac{1}{2} \int_X q \log(1 - 2e^{2\varphi}) dx - \frac{1}{2} \log 2.$$

□

Proposition 5. Suppose μ has density $p(x)$ and ν has density $q(x)$. The influence function for J_{NS} is

$$\Psi_{\text{NS}}(x) = \log \frac{p(x)}{q(x)}.$$

Proof. The result follows from Lemma 1:

$$\begin{aligned}
 & \left. \frac{d}{d\epsilon} J_{\text{NS}}(\mu + \epsilon\chi) \right|_{\epsilon=0} \\
 &= \left. \frac{d}{d\epsilon} \int_X (p + \epsilon\chi) \log \frac{p + \epsilon\chi}{q} dx \right|_{\epsilon=0} \\
 &= \int_X \left[\chi \log \frac{p}{q} + \chi \right] dx \\
 &= \int_X \left[\log \frac{p}{q} + 1 \right] d\chi \\
 &= \int_X \left[\log \frac{p}{q} \right] d\chi.
 \end{aligned}$$

□

Proposition 7. *The influence function for J_{W} is the Kantorovich potential corresponding to the optimal transport from μ to ν .*

Proof. See Santambrogio (2015), Proposition 7.17. □

Proposition 8. *The convex conjugate of J_{W} is*

$$J_{\text{W}}^*(\varphi) = \mathbb{E}_{x \sim \nu}[\varphi(x)] + \{ \|\varphi\|_L \leq 1 \}.$$

Proof. Using Kantorovich–Rubinstein duality, we have that

$$\begin{aligned}
 J_{\text{W}}(\mu) &= \sup_{\|\varphi\|_L \leq 1} \left[\int_X \varphi d\mu - \int_X \varphi d\nu \right] \\
 &= \sup_{\varphi} \left[\int_X \varphi d\mu - \int_X \varphi d\nu - \{ \|\varphi\|_L \leq 1 \} \right],
 \end{aligned}$$

where we use the notation

$$\{A\} = \begin{cases} 0 & A \text{ is true,} \\ \infty & A \text{ is false.} \end{cases}$$

By Lemma 2,

$$J_{\text{W}}^*(\varphi) = \int_X \varphi d\nu + \{ \|\varphi\|_L \leq 1 \}.$$

□

Proposition 10. *The influence function for J_{VI} is*

$$\Psi_{\text{VI}}(z) = \log \frac{q(z)}{p(x|z)p(z)}.$$

Proof. The result follows from Lemma 1:

$$\begin{aligned}
 & \left. \frac{d}{d\epsilon} J_{\text{VI}}(q + \epsilon\chi) \right|_{\epsilon=0} \\
 &= \left. \frac{d}{d\epsilon} \int (q(z) + \epsilon\chi(z)) \log \frac{q(z) + \epsilon\chi(z)}{p(z|x)} dz \right|_{\epsilon=0} \\
 &= \int \left[\chi(z) \log \frac{q(z) + \epsilon\chi(z)}{p(z|x)} + \chi(z) \right] dz \Big|_{\epsilon=0} \\
 &= \int \left[\log \frac{q(z)}{p(z|x)} + 1 \right] \chi(z) dz \\
 &= \int \left[\log \frac{q(z)}{p(x|z)p(z)} + \log p(x) + 1 \right] \chi(z) dz \\
 &= \int \log \frac{q(z)}{p(x|z)p(z)} \chi(z) dz.
 \end{aligned}$$

□

Proofs continue on the following page.

Proposition 13. *The influence function for J_{RL} is*

$$\Psi_{\text{RL}}(s, a) = -\frac{\sum_{t=0}^{\infty} \gamma^t p_t^\pi(s)}{\pi(s)} (Q^\pi(s, a) - V^\pi(s)),$$

where Q^π is the state-action value function, V^π is the state value function, and p_t^π is the marginal distribution of states after t steps, all under the policy π .

Proof. First, we note that

$$\begin{aligned} & \left. \frac{d}{d\epsilon} (\pi + \epsilon\chi)(a|s) \right|_{\epsilon=0} \\ &= \left. \frac{d}{d\epsilon} \frac{\pi(a, s) + \epsilon\chi(s, a)}{\pi(s) + \epsilon\chi(s)} \right|_{\epsilon=0} \\ &= \frac{\chi(s, a) - \chi(s)\pi(a|s)}{\pi(s)}, \end{aligned}$$

where we abuse notation to denote $\chi(s) = \int \chi(s, a') da'$.

We have

$$-J_{\text{RL}} = \mathbb{E} \left[\sum_{t=1}^{\infty} \gamma^{t-1} r_t \right],$$

or, plugging in the measure,

$$-J_{\text{RL}} = \int \sum_{t=1}^{\infty} \gamma^{t-1} r_t p_0(s_0) \prod_{j=1}^{\infty} p(s_j, r_j | s_{j-1}, a_j) \prod_{k=1}^{\infty} \pi(a_k | s_{k-1}).$$

The integral is over all free variables; we omit them here and in the following derivation for conciseness.

In computing $\left. \frac{d}{d\epsilon} J_{\text{RL}}(\pi + \epsilon\chi) \right|_{\epsilon=0}$, the product rule dictates that a term appear for every k , in which $\pi(a_k | s_{k-1})$ is replaced with $\left. \frac{d}{d\epsilon} (\pi + \epsilon\chi)(a_k | s_{k-1}) \right|_{\epsilon=0}$. Hence:

$$\begin{aligned} & \left. \frac{d}{d\epsilon} J_{\text{RL}}(\pi + \epsilon\chi) \right|_{\epsilon=0} \\ &= \int \sum_{t=1}^{\infty} \gamma^{t-1} r_t p_0(s_0) \prod_{j=1}^{\infty} p(s_j, r_j | s_{j-1}, a_j) \\ & \quad \times \sum_{k=1}^{\infty} \frac{\chi(s_{k-1}, a_k) - \chi(s_{k-1})\pi(a_k | s_{k-1})}{\pi(s_{k-1})} \prod_{\substack{\ell=1 \\ \ell \neq k}}^{\infty} \pi(a_\ell | s_{\ell-1}) \\ &= \sum_{k=1}^{\infty} \int \sum_{t=1}^{\infty} \gamma^{t-1} r_t p_0(s_0) \prod_{j=1}^{\infty} p(s_j, r_j | s_{j-1}, a_j) \\ & \quad \times \frac{\chi(s_{k-1}, a_k) - \chi(s_{k-1})\pi(a_k | s_{k-1})}{\pi(s_{k-1})} \prod_{\substack{\ell=1 \\ \ell \neq k}}^{\infty} \pi(a_\ell | s_{\ell-1}), \end{aligned}$$

reordering the summations. Note that for $t < k$, the summand vanishes:

$$\begin{aligned}
 & \int \prod_{j=k}^{\infty} p(s_j, r_j | s_{j-1}, a_j) \\
 & \quad \times (\chi(s_{k-1}, a_k) - \chi(s_{k-1})\pi(a_k | s_{k-1})) \prod_{\ell=k+1}^{\infty} \pi(a_\ell | s_{\ell-1}) \\
 & = \int (\chi(s_{k-1}, a_k) - \chi(s_{k-1})\pi(a_k | s_{k-1})) \\
 & = \int (\chi(s_{k-1}) - \chi(s_{k-1})) \\
 & = 0,
 \end{aligned}$$

since all the variables $a_k, r_k, s_k, a_{k+1}, r_{k+1}, s_{k+1}, \dots$ integrate away to 1. This yields:

$$\begin{aligned}
 & -\frac{d}{d\epsilon} J_{\text{RL}}(\pi + \epsilon\chi) \Big|_{\epsilon=0} \\
 & = \sum_{k=1}^{\infty} \int \sum_{t=k}^{\infty} \gamma^{t-1} r_t p_0(s_0) \prod_{j=1}^{\infty} p(s_j, r_j | s_{j-1}, a_j) \\
 & \quad \times \frac{\chi(s_{k-1}, a_k) - \chi(s_{k-1})\pi(a_k | s_{k-1})}{\pi(s_{k-1})} \prod_{\substack{\ell=1 \\ \ell \neq k}}^{\infty} \pi(a_\ell | s_{\ell-1}).
 \end{aligned}$$

Then, substituting the marginal distribution (note s_{k-1} is not integrated)

$$p_{k-1}^\pi(s_{k-1}) = \int \prod_{j=1}^{k-1} p(s_j, r_j | s_{j-1}, a_j) \prod_{\ell=1}^{k-1} \pi(a_\ell | s_{\ell-1}),$$

we obtain

$$\begin{aligned}
 & -\frac{d}{d\epsilon} J_{\text{RL}}(\pi + \epsilon\chi) \Big|_{\epsilon=0} \\
 & = \sum_{k=1}^{\infty} \int \sum_{t=k}^{\infty} \gamma^{t-1} r_t p_{k-1}^\pi(s_{k-1}) \prod_{j=k}^{\infty} p(s_j, r_j | s_{j-1}, a_j) \\
 & \quad \times \frac{\chi(s_{k-1}, a_k) - \chi(s_{k-1})\pi(a_k | s_{k-1})}{\pi(s_{k-1})} \prod_{\ell=k+1}^{\infty} \pi(a_\ell | s_{\ell-1}).
 \end{aligned}$$

Let us rename the integration variables by decreasing their indices by $k - 1$:

$$\begin{aligned}
 & -\frac{d}{d\epsilon} J_{\text{RL}}(\pi + \epsilon\chi) \Big|_{\epsilon=0} \\
 & = \sum_{k=1}^{\infty} \int \sum_{t=1}^{\infty} \gamma^{t+k-2} r_t p_{k-1}^\pi(s_0) \prod_{j=1}^{\infty} p(s_j, r_j | s_{j-1}, a_j) \\
 & \quad \times \frac{\chi(s_0, a_1) - \chi(s_0)\pi(a_1 | s_0)}{\pi(s_0)} \prod_{\ell=2}^{\infty} \pi(a_\ell | s_{\ell-1}).
 \end{aligned}$$

Substituting in

$$\begin{aligned}
 V^\pi(s_0) & = \int \sum_{t=1}^{\infty} \gamma^{t-1} r_t \prod_{j=1}^{\infty} p(s_j, r_j | s_{j-1}, a_j) \prod_{\ell=1}^{\infty} \pi(a_\ell | s_{\ell-1}), \\
 Q^\pi(s_0, a_1) & = \int \sum_{t=1}^{\infty} \gamma^{t-1} r_t \prod_{j=1}^{\infty} p(s_j, r_j | s_{j-1}, a_j) \prod_{\ell=2}^{\infty} \pi(a_\ell | s_{\ell-1}),
 \end{aligned}$$

we obtain

$$\begin{aligned} & -\frac{d}{d\epsilon} J_{\text{RL}}(\pi + \epsilon\chi) \Big|_{\epsilon=0} \\ & = \sum_{k=1}^{\infty} \int \gamma^{k-1} p_{k-1}^{\pi}(s_0) \frac{Q^{\pi}(s_0, a_1)\chi(s_0, a_1) - V^{\pi}(s_0)\chi(s_0)}{\pi(s_0)}. \end{aligned}$$

Finally, by [Lemma 1](#), we obtain that

$$\Psi_{\text{RL}}(s, a) = -\frac{\sum_{k=0}^{\infty} \gamma^k p_k^{\pi}(s)}{\pi(s)} (Q^{\pi}(s, a) - V^{\pi}(s)).$$

□

Proposition 16. *The convex conjugate of J_{RL} is*

$$J_{\text{RL}}^*(\varphi) = (1 - \gamma)\mathbb{E}_{p_0(s)} V_{\varphi}(s) + \{V_{\varphi} \text{ exists}\},$$

where V_{φ} is the unique solution to $\varphi = -\mathcal{A}V_{\varphi}$, if it exists.

Proof. As mentioned in the text, we set the arbitrary distribution $\pi(s) = (1 - \gamma) \sum_{t=0}^{\infty} \gamma^t p_t^{\pi}(s)$. In doing so, $\pi(s, a)$ becomes a state-action *occupancy measure* that describes the frequency of encounters of the state-action pair (s, a) over trajectories governed by the policy $\pi(a|s)$. It is known that there is a bijection between occupancy measures $\pi(s, a)$ and policies $\pi(a|s)$ ([Syed et al., 2008](#); [Ho & Ermon, 2016](#)).

We can enforce this setting by redefining

$$J_{\text{RL}}(\pi) = -\mathbb{E} \sum_{t=1}^{\infty} \gamma^{t-1} r_t + \left\{ \forall s : \pi(s) = (1 - \gamma) \sum_{t=0}^{\infty} \gamma^t p_t^{\pi}(s) \right\},$$

where again $\{\cdot\}$ is the convex indicator function. This equation can be rewritten as

$$J_{\text{RL}}(\pi) = -\mathbb{E}_{\pi(s,a)} R(s, a) + \left\{ \forall s' : \pi(s') = (1 - \gamma)p_0(s') + \gamma \mathbb{E}_{\pi(s,a)} p(s'|s, a) \right\},$$

where $R(s, a) = \mathbb{E}_{p(s',r|s,a)}[r]$. The constraint is known as the *Bellman flow equation*. This formulation is convex, as it is the sum of an affine function and an indicator of a convex set (indeed, an affine subspace).

We recall $-\varphi = \mathcal{A}V_{\varphi}$, where $\mathcal{A}V(s, a) = \mathbb{E}_{p(s',r|s,a)}[r + \gamma V(s')] - V(s)$. Now, V_{φ} is uniquely defined by φ if a solution to the equation exists. To see this, note that V_{φ} is the fixed point of the Bellman operator \mathcal{T}^a defined by

$$(\mathcal{T}^a V)(s) = (R + \varphi)(s, a) + \gamma \mathbb{E}_{p(s'|s,a)} V(s'),$$

which is contractive and therefore has a unique fixed point. A representation of V_{φ} may be obtained via fixed point iteration using \mathcal{T}^a for an arbitrary action a :

$$V_{\varphi}(s) = \lim_{k \rightarrow \infty} (\mathcal{T}^a)^k 0 = \mathbb{E}^a \sum_{t=1}^{\infty} \gamma^{t-1} (R + \varphi)(s_t, a),$$

where the expectation is taken under the deterministic policy a .

We rewrite J_{RL} using a Lagrange multiplier $V(s)$

$$\begin{aligned} J_{\text{RL}}(\pi) & = -\mathbb{E}_{\pi(s,a)} R(s, a) + \sup_V \int V(s') \left[\pi(s') - (1 - \gamma)p_0(s') - \gamma \mathbb{E}_{\pi(s,a)} p(s'|s, a) \right] ds' \\ & = \sup_V -\mathbb{E}_{\pi(s,a)} R(s, a) + \mathbb{E}_{\pi(s)} V(s) - (1 - \gamma)\mathbb{E}_{p_0(s)} V(s) - \gamma \mathbb{E}_{\pi(s,a)} \mathbb{E}_{p(s'|s,a)} V(s') \\ & = \sup_{\varphi} \mathbb{E}_{\pi(s,a)} \varphi(s, a) - (1 - \gamma)\mathbb{E}_{p_0(s)} V_{\varphi}(s) - \{V_{\varphi} \text{ exists}\}. \end{aligned}$$

Probability Functional Descent

Note that $(1 - \gamma)\mathbb{E}_{p_0(s)}V_\varphi(s) + \{V_\varphi \text{ exists}\}$ is convex in φ ; this stems from the fact that

$$V_{\alpha\varphi+(1-\alpha)\varphi'} = \alpha V_\varphi + (1 - \alpha)V_{\varphi'}.$$

The result follows from [Lemma 2](#).

□