

A Derivation of the Quadratic Form

Proof of Lemma 3. Consider function $h : \mathbb{R} \rightarrow \mathbb{R}$,

$$h(w) = c^\top \widehat{V}_{\bar{\pi}} w + v_{\bar{\pi},c}^\top (L_{\pi_w}(\widehat{V}_{\bar{\pi}} w) - \widehat{V}_{\bar{\pi}} w).$$

For a scalar w , define $\widehat{Q}_{\bar{\pi}}(x, a, w) = r(x) + \gamma w P_{(x,a)} \widehat{V}_{\bar{\pi}}$. Substituting for the Bellman operator L_{π_w} (see Section 1.1), we obtain

$$h(w) = c^\top \widehat{V}_{\bar{\pi}} w - v_{\bar{\pi},c}^\top \widehat{V}_{\bar{\pi}} w + \sum_x v_{\bar{\pi},c}(x) \sum_a \nu(a|x) \left(1 + \widehat{Q}_{\bar{\pi}}(x, a, w) - \mathbf{E}_{\nu(\cdot|x)} \widehat{Q}_{\bar{\pi}}(x, \cdot, w)\right) \widehat{Q}_{\bar{\pi}}(x, a, w).$$

Because $\widehat{Q}_{\bar{\pi}}(x, a, w) = r(x) + \gamma w P_{(x,a)} \widehat{V}_{\bar{\pi}}$, h is quadratic in w , so we can write it as $h(w) = (1/2)w^\top B w + g^\top w + f$ for some choice of parameters B , g , and f . We have that

$$\begin{aligned} \mathbf{E}_{\nu(\cdot|x)} \widehat{Q}_{\bar{\pi}}(x, \cdot, w) &= \sum_a \nu(a|x) \widehat{Q}_{\bar{\pi}}(x, a, w) \\ &= \sum_a \nu(a|x) (r(x) + \gamma w P_{(x,a)} \widehat{V}_{\bar{\pi}}) \\ &= r(x) + \gamma w \mathbf{E}_{\nu(\cdot|x)} (P \widehat{V}_{\bar{\pi}}). \end{aligned}$$

Also, we have

$$\begin{aligned} \mathbf{E}_{\nu(\cdot|x)} \widehat{Q}_{\bar{\pi}}^2(x, \cdot, w) &= \sum_a \nu(a|x) (\widehat{Q}_{\bar{\pi}}(x, a, w))^2 \\ &= \sum_a \nu(a|x) (r(x) + \gamma w P_{(x,a)} \widehat{V}_{\bar{\pi}})^2 \\ &= \sum_a \nu(a|x) \left(r(x)^2 + \gamma^2 w^2 (P_{(x,a)} \widehat{V}_{\bar{\pi}})^2 + 2\gamma w r(x) P_{(x,a)} \widehat{V}_{\bar{\pi}} \right) \\ &= r(x)^2 + 2\gamma w r(x) \mathbf{E}_{\nu(\cdot|x)} (P \widehat{V}_{\bar{\pi}}) + \gamma^2 w^2 \mathbf{E}_{\nu(\cdot|x)} (P \widehat{V}_{\bar{\pi}})^2. \end{aligned}$$

Thus,

$$\begin{aligned} \mathbf{Var}_{\nu(\cdot|x)} \widehat{Q}_{\bar{\pi}}(x, \cdot, w) &= \mathbf{E}_{\nu(\cdot|x)} \widehat{Q}_{\bar{\pi}}^2(x, \cdot, w) - (\mathbf{E}_{\nu(\cdot|x)} \widehat{Q}_{\bar{\pi}}(x, \cdot, w))^2 \\ &= r(x)^2 + 2\gamma w r(x) \mathbf{E}_{\nu(\cdot|x)} (P \widehat{V}_{\bar{\pi}}) + \gamma^2 w^2 \mathbf{E}_{\nu(\cdot|x)} (P \widehat{V}_{\bar{\pi}})^2 \\ &\quad - r(x)^2 - \gamma^2 w^2 (\mathbf{E}_{\nu(\cdot|x)} (P \widehat{V}_{\bar{\pi}}))^2 - 2\gamma w r(x) \mathbf{E}_{\nu(\cdot|x)} (P \widehat{V}_{\bar{\pi}}) \\ &= \gamma^2 w^2 \mathbf{Var}_{\nu(\cdot|x)} (P \widehat{V}_{\bar{\pi}}). \end{aligned}$$

Further, we have that

$$\begin{aligned} h(w) - c^\top \widehat{V}_{\bar{\pi}} w + v_{\bar{\pi},c}^\top \widehat{V}_{\bar{\pi}} w &= \sum_x v_{\bar{\pi},c}(x) \sum_a \nu(a|x) \left(1 + \widehat{Q}_{\bar{\pi}}(x, a, w) - \mathbf{E}_{\nu(\cdot|x)} \widehat{Q}_{\bar{\pi}}(x, \cdot, w)\right) \widehat{Q}_{\bar{\pi}}(x, a, w) \\ &= \sum_x v_{\bar{\pi},c}(x) \sum_a \nu(a|x) \left(\widehat{Q}_{\bar{\pi}}(x, a, w) + (\widehat{Q}_{\bar{\pi}}(x, a, w))^2 - \widehat{Q}_{\bar{\pi}}(x, a, w) \mathbf{E}_{\nu(\cdot|x)} \widehat{Q}_{\bar{\pi}}(x, \cdot, w) \right) \\ &= \sum_x v_{\bar{\pi},c}(x) \left(\mathbf{E}_{\nu(\cdot|x)} \widehat{Q}_{\bar{\pi}}(x, \cdot, w) + \mathbf{E}_{\nu(\cdot|x)} \widehat{Q}_{\bar{\pi}}(x, \cdot, w)^2 - (\mathbf{E}_{\nu(\cdot|x)} \widehat{Q}_{\bar{\pi}}(x, \cdot, w))^2 \right) \\ &= \sum_x v_{\bar{\pi},c}(x) \mathbf{E}_{\nu(\cdot|x)} \widehat{Q}_{\bar{\pi}}(x, \cdot, w) + \sum_x v_{\bar{\pi},c}(x) \mathbf{Var}_{\nu(\cdot|x)} \widehat{Q}_{\bar{\pi}}(x, \cdot, w), \end{aligned}$$

and therefore

$$h(w) = c^\top \widehat{V}_{\bar{\pi}} w + \sum_x v_{\bar{\pi},c}(x) \mathbf{E}_{\nu(\cdot|x)} \widehat{Q}_{\bar{\pi}}(x, \cdot, w) + \sum_x v_{\bar{\pi},c}(x) \mathbf{Var}_{\nu(\cdot|x)} \widehat{Q}_{\bar{\pi}}(x, \cdot, w) - v_{\bar{\pi},c}^\top \widehat{V}_{\bar{\pi}} w,$$

or alternatively,

$$h(w) = v_{\bar{\pi},c}^\top r + (c^\top \widehat{V}_{\bar{\pi}} - v_{\bar{\pi},c}^\top \widehat{V}_{\bar{\pi}} + \gamma \mathbf{E}_{v_{\bar{\pi},c}} (P \widehat{V}_{\bar{\pi}})) w + w^2 \sum_x v_{\bar{\pi},c}(x) \mathbf{Var}_{\nu(\cdot|x)} \widehat{Q}_{\bar{\pi}}(x, \cdot).$$

We therefore obtain

$$\begin{aligned}
 f &= v_{\hat{\pi},c}^\top r, \\
 g &= c^\top \hat{V}_{\hat{\pi}} - v_{\hat{\pi},c}^\top \hat{V}_{\hat{\pi}} + \gamma \mathbf{E}_{v_{\hat{\pi},c}}(P^\nu \hat{V}_{\hat{\pi}}), \\
 B &= 2 \sum_x v_{\hat{\pi},c}(x) \mathbf{Var}_{\nu(\cdot|x)} \hat{Q}_{\hat{\pi}}(x, \cdot).
 \end{aligned}$$

□