## Supplement

Here we provide the derivations that were omitted in the main text. The following two identities will be used repeatedly. Since $\mu$ is stationary, we have

$$\sum_x \mu(x, \mathbf{w}) \pi(x'|x, \mathbf{w}) = \mu(x', \mathbf{w}) \tag{42}$$

Since $\sum_{x'} \pi(x'|x, \mathbf{w}) = 0$ for all $\mathbf{w}$, we can differentiate and obtain

$$\sum_{x'} \nabla_{\mathbf{w}} \pi(x'|x, \mathbf{w}) = 0 \tag{43}$$

Again, we will suppress the dependence on $\mathbf{w}$.

**Proof of Theorem 1:**

Differentiating the Bellman equation (2) yields

$$\nabla_{\mathbf{w}} c + \nabla_{\mathbf{w}} v(x) = \sum_{x'} \nabla_{\mathbf{w}} \pi(x'|x) \left( \log \frac{\pi(x'|x)}{p(x'|x)} + v(x') \right) \tag{44}$$

$$+ \pi(x'|x) \left( \frac{\nabla_{\mathbf{w}} \pi(x'|x)}{\pi(x'|x)} + \nabla_{\mathbf{w}} v(x') \right)$$

$$= \sum_{x'} \nabla_{\mathbf{w}} \pi(x'|x) \left( \log \frac{\pi(x'|x)}{p(x'|x)} + v(x') \right) + \pi(x'|x, \mathbf{w}) \nabla_{\mathbf{w}} v(x')$$

To obtain the last equation we used (43). Now we move $\nabla_{\mathbf{w}} v(x)$ on the right side of (44), multiply by $\mu(x)$ and sum over $x$. Noting that

$$\sum_{x,x'} \mu(x) \pi(x'|x) \nabla_{\mathbf{w}} v(x') = \sum_x \mu(x) \nabla_{\mathbf{w}} v(x) \tag{45}$$

which follows from (42), the LMDP policy gradient is as given in (5).

**Proof of Theorem 2:**

Using the identity $\nabla_{\mathbf{w}} \pi = \pi \nabla_{\mathbf{w}} \log \pi$, equation (5) can also be written as

$$\nabla_{\mathbf{w}} c = \sum_{x,x'} \mu(x, x') \nabla_{\mathbf{w}} \log \pi(x'|x) \left( \log \frac{\pi(x'|x)}{p(x'|x)} + v(x') \right) \tag{46}$$

With the policy parameterization (7), we have

$$\nabla_{\mathbf{w}} \log \pi(x'|x) = \nabla_{\mathbf{w}} \left( \log p(x'|x) - \mathbf{w}^{\mathsf{T}} \mathbf{f}(x') - \log \sum_y p(y|x) \exp\left(-\mathbf{w}^{\mathsf{T}} \mathbf{f}(y)\right) \right) \tag{47}$$

$$= -\mathbf{f}(x') + \sum_y \frac{p(y|x) \exp\left(-\mathbf{w}^{\mathsf{T}} \mathbf{f}(y)\right)}{\sum_s p(s|x) \exp\left(-\mathbf{w}^{\mathsf{T}} \mathbf{f}(s)\right)} \mathbf{f}(y)$$

$$= \Pi[\mathbf{f}](x) - \mathbf{f}(x')$$

Substituting (47) in (46) and using the fact that

$$\sum_x \mu(x) \log \left( \sum_y p(y|x) \exp\left(-\mathbf{w}^{\mathsf{T}} \mathbf{f}(y)\right) \right) \sum_{x'} \nabla_{\mathbf{w}} \pi(x'|x) = 0 \tag{48}$$

which follows from (43), the gradient is as given in (9).

**Proof of Theorem 4:**

Using (47), equation (17) can be written as

$$\nabla_{\mathbf{w}} c = \sum_{x,x'} \mu(x, x') \nabla_{\mathbf{w}} \log \pi(x'|x) \left( \nabla_{\mathbf{w}} \log \pi(x'|x) - \Pi[\mathbf{f}](x) \right)^{\mathsf{T}} (\mathbf{w} - \mathbf{r}) \tag{49}$$

$$= G(\mathbf{w})(\mathbf{w} - \mathbf{r}) - \sum_{x,x'} \mu^{\pi}(x) \nabla_{\mathbf{w}} \pi(x'|x) \Pi[\mathbf{f}](x)^{\mathsf{T}} (\mathbf{w} - \mathbf{r})$$

The second term is zero because of (43), thus we have equation (19).