Our Q-function:

$$Q(s,a) = r(s,a) + \beta \sum_{s'} P(s'|s,a)V(s'),$$

where

$$V(s) = E\{\max_{a \in A} Q(s,a) + \epsilon_a\}.$$

Thus,

$$V(s) = \sum_a P(a|s)E[Q(s,a) + \epsilon_a|Y^* = a]$$

$$= r^\pi(s) + \sum_a E[\epsilon_a|Y^* = a] + \beta \sum_{s'} P^\pi(s'|s)V(s')$$

$$= r^\pi(s) + |A|\gamma - \sum_a \log(\pi(a|s)) + \beta \sum_{s'} P^\pi(s'|s)V(s')$$

After normalization, what we really have is:

$$V(s) = r^\pi(s) - \sum_a \log(\pi(a|s)) + \beta \sum_{s'} P^\pi(s'|s)V(s').$$

For any given $\pi$, define $H^\pi : V \to V$, $V : S \to R$, we have

$$[H^\pi v](s) = r^\pi(s) - \sum_a \log(\pi(a|s)) + \beta \sum_{s'} P^\pi(s'|s)v(s'),$$

which is a contraction mapping. For each $\theta$ in $r_\theta$, IRL inner loop is trying to find $\pi^*$, a randomized Markovian policy, such that

$$H^{\pi^*} v^* = v^*$$

$$\pi^*(a|s) \propto e^{r(s,a) + \beta E_{s'|s} v^*(s')}.$$

This may be viewed as a perturbed version of MDP, where we have $\sup_\pi H^\pi v^* = v^*$ in our textbooks. To identify $\theta$, they essentially solved

$$\min_\theta KL(\pi_\theta^* || \hat{\pi})$$

$$s.t. H^{\pi^*} v^* = v^*$$

$$\pi^*(a|s) \propto e^{r(s,a) + \beta E_{s'|s,a} v^*(s')}.$$

The inner loop is expensive. Enforce $\pi_\theta^* = \hat{\pi}$, we have

$$H_\theta^{\hat{\pi}} v_\theta^{\hat{\pi}} = v_\theta^{\hat{\pi}},$$

which is a set of system linear equations depending on $\theta$ (note this is not policy iteration; this is only policy evaluation step). Since $[I - \beta P^{\hat{\pi}}]^{-1}$ only needs to be inverted once, $v_\theta^{\hat{\pi}}$ can be obtained efficiently for different $\theta$ (a simple matrix * vector). This is why CCP method is much faster than NFXP. Accordingly, $\pi_\theta(a|s) \propto e^{r_\theta(s,a) + \beta E_{s'|s} v_\theta^{\hat{\pi}}(s')}$, you now select $\theta$ to maximize log-likelihood. That is,

$$\max_\theta \sum_{t=1}^T \log \pi_\theta(a_t|s_t)$$

$$s.t. H_\theta^{\hat{\pi}} v_\theta^{\hat{\pi}} = v_\theta^{\hat{\pi}}$$

$$\pi_\theta(a|s) \propto e^{r_\theta(s,a) + \beta E_{s'|s,a} v_\theta^{\hat{\pi}}(s')}$$

Given $v$, construct $\pi$ is really easy. The real bottlenecks are (1) taking gradient of $\log \pi_\theta(a_t|s_t)$ and (2) inverting $H_\theta^{\hat{\pi}} v_\theta^{\hat{\pi}} = v_\theta^{\hat{\pi}}$.

1

For (1), the reason is $\pi_\theta$ being implicit, while for (2) inversion is expensive for large or infinite $S$. Recall Gaussian elimination is $O(n^3)$ (or maybe 2.6 or 2.8, i forget the exact number). There are some ways for efficiently inverting matrix, which mainly hinge upon utilizing special structure of a matrix, such as sparsity, symmetry, etc. Actor-critic assumes $\pi$ is in the form of $e^{\omega^T\Phi}$, i.e., parameterize $Q$ function. For (2), two popular ways in the literature, neural network or linear approximation. I am an idiot on neural network, so let's proceed with linear approximation. In fact, this is a main idea of approximate dynamic programming (ADP) and I highly suspect many algorithmic and theoretical results from ADP can be extended to this case with a slightly modified $H$ operator.

It will be great if we know some features of the reward function $f$. Remember, successful implementation of ADP requires you have a good understanding of your problem's structural results, such as convexity, piecewise linearity, etc. Our belief-based IRL model is indeed convex and we haven't examined how to utilize it.

Build Krylov space $\{f-\sum_a \log(\pi(a|S)), \beta P^{\hat\pi}(f-\sum_a \log(\pi(a|S))), ..., \beta^K P^{\hat\pi,K-1}(f-\sum_a \log(\pi(a|S)))\} = \{\Phi_1, ..., \Phi_K\}$.

$\forall \theta$, we want to find $\tilde{v}_\theta^{\hat\pi} = \omega^T\Phi \approx v_\theta^{\hat\pi}$, where $\omega \in R^K, K << |S|$. Namely,

$$\tilde{v}_\theta^{\hat\pi} = \arg \max_{u \in Krylov} ||u - v_\theta^{\hat\pi}||.$$

Equivalently, we seek $\omega^*$ such that

$$\omega^* \in \arg\min_\omega ||\omega^T\Phi(S) - (r_\theta(S) - \sum_a \log(\hat\pi(a|S)) + \beta \int_S P^{\hat\pi}(s'|s)\omega^T\Phi(s'))ds'||_D$$

$$\arg\min_\omega ||\omega^T(\Phi - \beta P^{\hat\pi}\Phi) - (r_\theta - \sum_a \log(\hat\pi(a|S))||_D$$

where $D$ has the $d_s$ on its diagonal and $d = (d_s)$ is the stationary distribution of $P^{\hat\pi}$.

Note while $|S|$ can be infinite, we only need $K$ equations to uniquely determine $\omega$, and taking the inverse of $(\Phi - \beta P^{\hat\pi}\Phi)$ is manageable since $K$ is small and we only need to do it once.

Now we construct

$$\pi_\theta(a|s) = \frac{e^{r_\theta(s,a)+\beta\omega^T P^a\Phi}}{\sum_{a'} e^{r_\theta(s,a')+\beta\omega^T P^{a'}\Phi}}$$

Derivative of softmax has an analytical expression.. So update $\theta$, continue...

If $r_\theta \in Krylov$ space we construct, we will have

$$\pi_\theta(a|s) = \frac{e^{(\theta+\beta\omega^T P^a)\Phi}}{\sum_{a'} e^{(\theta+\beta\omega^T P^{a'})\Phi}},$$

which will allow us to directly use actor-critic to calculate $v$. We can test which way is faster.