# IV. Online ADP Algorithm

① TD Error
② Value Fcn Approx.
$\left.\right\} \rightarrow$ Policy Itr.
$\hookrightarrow$ Online ADP

At each time step $k$, collect data $(x_k, x_{k+1},$
$c(x_k, \pi(x_k)))$. Consider value fcn approx:
$V_\pi(x) = W^T \phi(x)$. Then the TD err is
$$e_k = c(x_k, \pi(x_k)) + \gamma \cdot W^T \phi(x_{k+1}) - W^T \phi(x_k)$$

corresponds to linear regression model:

$$\Rightarrow c(x_k, \pi(x_k)) = [\phi(x_k) - \gamma \cdot \phi(x_{k+1})]^T W$$

re-arranged Bellman eqn: $V_k = c + \gamma \cdot V_{k+1}$

We can now write our first online RL algo that performs policy eval via supervised learning

## Online Policy Iter.

0) Initialization: Select an admissible control policy $\pi^0$. Set $m = 0$.

1) Policy Eval: Run control policy $\pi^m$ on environ/ system for one episode. Collect $L$ measured data tuples $(x_k, x_{k+1}, c(x_k, \pi^m(x_k)))$. Find least squares solution w.r.t. $W_m$ for regression model (a.k.a Bellman Eqn)

$$
\underbrace{\begin{bmatrix} \vdots \\ c(x_k, \pi^m(x_k)) \\ \vdots \end{bmatrix}}_{L \times 1} \overset{=C}{=} \underbrace{\begin{bmatrix} \vdots \\ [\phi(x_L) - \gamma \cdot \phi(x_{k+1})] \\ \vdots \end{bmatrix}^T}_{L \times n_w} \overset{=\Phi}{} \underbrace{W}_{n_w \times 1}
$$

$L > n_w$

written compactly as $C = \underset{\varepsilon}{\Phi} W$

For example, you can perform ordinary lsq.

$$W_m \leftarrow W^* = [\Phi^T \Phi]^{-1} \Phi^T G$$

2) Policy Improve: Find an improved policy via

$$\Pi^{m+1} = \arg\min_{\Pi(\cdot)} \{ c(x_k, \Pi(x_k)) + \gamma \cdot W_m^T \phi(x_{k+1}) \}$$

where $x_{k+1} = f(x_k, \Pi(x_k)) \quad \forall x_k \in X$

Set $m \leftarrow m+1$. Go to Step 1.

Rem: Besides OLS, you can also use recursive lsq, gradient method, ridge regression, LASSO regression.

Rem: In online ADP, the regressor

$$\left[\phi(x_k) - \gamma \cdot \phi(x_{k+1})\right] \text{ must be}$$

"persistently excited" for a soln to exist for lsq. This is a sufficient condition for $\Phi^T \Phi$ to be invertible.

Rem: Observe that Step 1 Policy Eval is model-free. We only require data $(x_k, x_{k+1}, c(\cdot, \cdot))$

However, Step 2 Policy Improvement is <u>NOT</u> model-free. We are required to solve:

$$\frac{\partial c}{\partial u}(x_k, \pi(x_k)) + \gamma \cdot W^{TD} \frac{\partial \phi}{\partial x}(x_{k+1}) \cdot \frac{\partial f}{\partial u}(x_k, \pi(x_k)) = 0$$

which requires knowledge of $c(\cdot, \cdot), f(\cdot, \cdot)$

<u>Rem</u> Step 2 still requires minimization for all $x_k \in X$. So we have only partially avoided the curse of dimensionality

This motivates fcn approx. for the control policy fcn $\pi(\cdot)$.

Called "actor neural net" by Werbos & Bertsekas.