# Reinforcement Learning
## Week 6 — Wednesday

## Error Projected Bellman

① doesn't involve true value, thus more stable than $E_{ms}$.

Scalar $V(s) = r(s) - g + E[V(s')]$

$\hookrightarrow$ with deterministic policy

matrix $\vec{V} = \vec{r} - g\vec{1} + P\cdot\vec{V}$
(vector)

$\swarrow$

One-step state distribution
Size: $|S|$ by $|S|$

$\vec{V} = B[\vec{V}]$

$\hookrightarrow$ Bellman (evaluator) operation

$$e_{PB} \overset{def}{=} \| \hat{v}(w) - \mathbb{P}\mathbb{B}\,\hat{v}(w) \|^2_{\rho*}$$

$$\hookrightarrow \text{norm vector}$$

$$= (\Delta v)^T D_{\rho*} \Delta v$$

where

~~$\Delta v = \hat{V}(w) - PB v$~~

$$= \sum_{s \in S} \rho^*(s) \cdot (\Delta v(s))^2$$

$\Delta v = \hat{V}(w) - PB\hat{V}w$

$$\doteq \underset{\rho^*}{\mathbb{E}}\left[\Delta v(s)\right]^2$$

Exceed expectations, derive P !

how to get $w^*$

① take the gradient of $\varepsilon_{PB}$ $\underbrace{\qquad\qquad\qquad}_{\nabla \varepsilon_{PB} \text{ of } (w)}$ $w$

② set the grad to zero because in local minima , gradient is $0$



$w^*$

. Then we can show that

$$w^* = X^{-1} y \qquad \longrightarrow X =$$

$$\sum_{s} p^*(s) \sum_{s^t} P(s'|s) \left[ f(s) \{ f(s) - f(s') \} \right]^T$$

$$\longrightarrow = F^T D_{p^*} (I - P) F$$

$$y = \sum_{S} p^*(s) \left[ (r(s) - g\mathbb{1}) \, f(s) \right]$$

$$= F^T \, D_{p^*} \, (r - g\mathbb{1})$$

## LSTD (Least Square Temporal Diff)

↳ Predict $\hat{w}^*$

$$\hat{w} = \boxed{\hat{X}^{-1} \, \hat{y}}$$

$$\hat{X} = \frac{1}{n} \overset{\sum_{i=1}^{n}}{F(s_i)} \left[ f(s_i) - f(s_i') \right]^T$$

$$\hat{y} = \frac{1}{n} \sum_{i=1}^{n} (r(s_i) - g\mathbb{1}) \cdot f(s_i)$$

$$\boxed{\hat{w}^* = \hat{X}^{-1} \cdot \hat{y}}$$