# RL (10/3/2025).

Agent $\Rightarrow$ Porciovo
$+$
Act

Environment

## Major in RL (schemes).

Policy Iteration     Value Iteration $---$ Both are DP setting in

↓ Derivative           ↓ Derivative

GPI (Generalized Policy Iteration)

inexact! (approximate).

$Q^*$ - learning.
$Q_X$ → optimality criteria : $g, b, \partial$ (total.)

$1^{st}$ optimality criteria.

In avg-rew setting, we don't care about $Q_g$. Thus we just use $Q_b$

Which is more intuitive?

Remember "no free lunch theorem", so sometimes PI is better, so is VI.

Back to inexact policy evaluation. $\longrightarrow$ value.

$\hookrightarrow$ It uses function approximator (like in supervised learning).

characteristic) Parametric              non-parametric.

$\hookrightarrow$ has fixed & finite number of params.

our focus    classified into:

linear        non-linear; e.g. NN.

notes: NN can be linear too!

What's being approximated : Value

$$\hat{V}(s;w) = w^T f(s) \approx v(s) \rightarrow \text{In matrix form:}$$

$$\vec{V}(w) = F \cdot w.$$

parameter vector

state-feature vector function.

$w \in W = \mathbb{R}^{dim(w)}$

$f : S \rightarrow \mathbb{R}^n.$

$n = dim(w)$ in linear setting.

$$\begin{bmatrix} v(s^0 w) \\ v(s' w) \\ \vdots \end{bmatrix} = \begin{bmatrix} f(s)^T \\ f(s')^T \\ \vdots \end{bmatrix} \begin{bmatrix} w_0 \\ w_1 \\ \vdots \\ w_{dim(w)-1} \end{bmatrix}$$

What's next after we define $\hat{v}$ ?

→ learn w ! → we need
   (training).         to define
                       error function,
                       of course to minimize
                       the error.

Yields:
optimal $w^* = \underset{w \in R^d}{\text{argmin}}\ \ell_x(w)$.

$\underbrace{\phantom{xxxxx}}$  ← error
                    for this week

Update
via SGD, ←    $\ell_{ms}$          $\ell_{PB}$
              (mean           (Projected
              square)          Bellman).

① $\ell_{ms}$ : (weighted)
              (we weight states with $\not{P}^*(s)$).

we need it bcs some states
may not be visited !

Stationary state dist
$\not{*} \neq *$
stationary    optimal

likelihood the agent
stays in S

$\ell_{ms}(w) = \sum_{s \in S} p^*(s) [v(s) - \hat{v}(s;w)]^2 = \mathbb{E}[v(S) - \hat{v}(S;w)]^2$
$\qquad\qquad\qquad\qquad\qquad\qquad\qquad\qquad\qquad S \sim p^*$

we use $p^*(s)$
instead of uniform weighting, i.e ~~~~~
multiply by $1/|S|$

For finding exact optimal policy:

~~we need~~

we may need weighted vector-norm:

$$\| \vec{v} - \vec{\hat{v}}(w) \|^2_{P^*}$$

$$= \left[ v - \hat{v}(w) \right]^T \cdot \underline{D_{P^*}} \left[ v - \hat{v}(w) \right].$$

↳ A Diagonal Square Matrix.

e.g. $\begin{bmatrix} P^*(s) & 0 & 0 & \cdots \\ 0 & P^*(s) & & \vdots \\ 0 & & \ddots & \vdots \\ & & & \cdots \end{bmatrix}$

During Learning (Training):

$\hookrightarrow$ update $\hat{w}$, using gradient <u>descent</u>.

$$\hat{w} \leftarrow \hat{w} - \alpha \cdot \nabla_{w} \ell_{ms}$$

where

hence ~~later~~
we use minus.
if ascent → plus.

$$\nabla \ell_{ms} = \nabla \underset{p\star}{\mathbb{E}} \left[ \left( v(s) - \hat{v}(s;w) \right)^2 \right]$$

$\quad$ exchange the
$\quad$ order of $\nabla$ vs $\mathbb{E}$

$$= \underset{p\star}{\mathbb{E}} \left[ \nabla \left( v(s) - \hat{v}(s;w) \right)^2 \right].$$

$$= \underset{p\star}{\mathbb{E}} \left[ 2 \cdot \left( v(s) - \hat{v}(s;w) \right) \cdot \left( -\nabla_{w} \hat{v}(s;w) \right) \right]$$

$\downarrow$ Introduce $\left( -\frac{1}{2} \right)$ as multiplier
in both sides!

$$-\frac{1}{2} \nabla \ell_{ms} : \underset{p\star}{\mathbb{E}} \left[ \underbrace{(v(s)}_{true} - \underbrace{\hat{v}(s;w))}_{approx.} \cdot \underbrace{\nabla_{w} \hat{v}(s;w)}_{} \right]$$

$\qquad \downarrow$

$\qquad = f(s)$ in linear
$\qquad\qquad$ case!

C.f. SL, ~~there~~

thus approximate
$\quad$ true value, e.g. with <u>TD</u>,

① there are no training
$\qquad$ data, so $\boxed{v(s) \text{ is unknown}}$ ! $\underline{\text{MC, etc.}}$

② in RL, data isn't i.i.d, i.e. <u>Markovian</u>,
$\qquad$ and collected by the <u>agent</u> itself!

For ①; How to approx. the true value?

$$\underset{\approx}{\overset{\curvearrowright}{z}} \ \underset{\approx}{\mathbb{E}} \left[ \left\{ \left\{ r(s,A) - \hat{g} + \hat{v}(s';w) \right\} - \hat{v}(s;w) \right\} \cdot f(s) \right]$$

approx operator ~
operator ^
as we use $\hat{g}, \hat{v}$

use sampling of $(s,a,r',s')$.

approx. $v(s)$ based on BEE $\Rightarrow$ involve $\hat{v}(w)$
hence bootstrap

$$\underset{\approx}{\overset{\curvearrowright}{z}} \ \left\{ \left\{ r(s,a) - \hat{g} + \hat{v}(s';w) - \hat{v}(s;w) \right\} f(s) \right\}$$

Sampling approximation.

term
4) bootstrap = using estimation ~~to esti~~ of
a value to estimate
another one -

# RL - Week 6 (12-3-2025).

## Implementations in RL:

1.) SL : approx policy evaluation ← caveats
   - ① not i.i.d.
   - ② true val unknown

2.) UL :- clustering states (aggregating, e.g. Jabodetabek)

   - State representation/feature extraction.

---

Continuing on error : Projected Bellman Error

~ In contrast with $\ell_{ms}$, $\ell_{PB}$ doesn't involve true value → more stable.

## Bellman Operator.

From BEE :  → correlating value of current with next state.

$$v(s) = r(s,a) - g + \mathbb{E}[v(s')]$$

$$\frac{TD}{I} = RHS - LHS = 0 \; ... \; theoritically...$$

TEMPOral Difference, i.e. $(t+1) - (t)$. → singular!

matrix (vector) form :

$$\vec{v} = \vec{r} - g\mathbf{1} + P\vec{v} \iff (I-P)v = r - g\mathbf{1}$$

↳ 1-step transition matrix

$$\vec{v} = \mathbb{B}[\vec{v}]$$

$\hookrightarrow$ Bellman (evaluator) operaton. on $\vec{v}$

$$e_{PB} = \left\| \hat{v}(w) - \underset{\uparrow}{P} \cdot \mathbb{B} \, \hat{v}(w) \right\|_{p^*}^2 \quad \text{--- norm vector.}$$

$\underbrace{\qquad\qquad\qquad\qquad\qquad}_{\Delta v}$  Projection operator.

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

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

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

Let $\vec{w} = [w_1 \; w_2]^T$ , $w_1, w_2 \in \mathbb{R}$.

Let $w = [w_1, w_2]^T$; $w_1, w_2 \in \mathbb{R}$.

$\mathbb{B}[\hat{v}(w)] \to$ is not in paramter space! $\to$ params $w$ can't represent $\mathbb{B}[\hat{v}_w]$



$\mathbb{B}[\hat{v}(w)]$

$w_1$

$\hat{v}(w)$

$\mathbb{P}\mathbb{B}\hat{v}$

$w_2$

to bring back $\hat{v}(w)$ after $\mathbb{B}[\hat{v}(w)]$ !

The projection operator $\mathbb{P}$ satisfies:

$$\mathbb{P}(v) = F \tilde{w}$$

true val? feature matrix.

later on will be cancel out!

where:

$$\tilde{w} = \underset{w \in W}{\arg\min} \left\{ \| F \cdot w - v \|^2_{p\star} \right\}.$$

$\hat{v}^k$

we can show that:

later on, the

$$P = F \left( F^T D_{p\star} F \right)^{-1} F^T D_{p\star}$$

Now, we should be able to get :

$$w^{\bigstar} = \underset{w}{\text{argmin}} \; \ell_{p_B}(w)_\star$$

↓ How?

① Take grad of $\ell_{PB}$

② Set the grad to 0 as in local minima, grad is 0.

We can find $w^{\bigstar}$ with closed form solution

$$w^{\bigstar} = X^{-1} \cdot y \rightarrow \text{Note}: \text{in DP setting.}$$

where

$$X = \sum_s p^{\bigstar}(s) \cdot \sum_{s'} p(s'|s) \left[ f(s) \{ f(s) - f(s') \}^T \right]$$

$$= \sum_s p^{\bigstar}(s) f(s) \left[ \left( f(s) - \sum_{s'} p(s'|s) f(s') \right)^T \right]$$

$$= F^T D_{p^{\bigstar}} (\mathbf{1} - P) \cdot F$$

and

$$y = \sum_s p^{\bigstar}(s) \left[ (r(s) - g\mathbf{1})^{an} \cdot f(s) \right]$$

$$= F^T D_{p^{\bigstar}} (r - g\mathbf{1})$$

In RL settings, we can't find $w^*$ using $X^{-1}y$ as before. Thus, we'll use LSTD, which is the approximation version of $X^{-1}y$.

Least Square Temporal Diff.

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

available because $X$ and $y$ are sampling-friendly, e.g. using expectation, i.e. ~~p~~ weight with ~~1~~.

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

$$\hat{y} = \frac{1}{n}\sum (r(s_i) - g) \cdot f(s_i).$$

Note: $n$ is number of samples. ~~last~~

Based on sample mean,

$$\hat{M} = \frac{1}{n}\sum_{i=1}^{n} X_i$$

which is unbiased!

Problem with LSTD:

$\hat{w} = \hat{X}^{-1} \hat{y}$, but sometimes $\hat{X}$ is not invertible.

Hence introduce $\varepsilon$ to add into it!

$$\hat{X} \leftarrow \hat{X} + \varepsilon \mathbf{I}.$$

Are there training & testing in RL?

In RL, there are no separation between ~~training~~ vs testing!

~~training~~: before convergence

after convergence.

If we really want to distinguish: Basic setting: In the same env

~~Even the~~ ~~environment for traing~~

~~But~~ maybe not in more challenging setting,

e.g. ① Non-stationary env

② Different but similar env.

Is tabular setting a special case of parametric function approximator?

Yes, where $f(s)$ uses one-hot encoding.
and $\dim(w) = |S|$ e.g.:

|  | $S = s^0$? | $S = s^1$? | . . . . |
|---|---|---|---|
| $f(s^0)$ | 1 | 0 | 0 . . . . |
| $f(s^1)$ | 0 | 1 | 0 . . . . |