During learning (training): we update $\hat{w}^n$ as
(online estimate)

No gradient descent method

$$\hat{w} \leftarrow \hat{w} - \alpha \nabla e_{ms}(w) \quad \text{here } mse$$

grad of error

Where the grad is $\quad \nabla \overset{def}{=} \frac{d}{dw}$

$$\nabla \underset{p^0}{E} \left[ V(s) - \hat{v}(s;w)^2 \right] = E\left[ \nabla (V(s) - \hat{v}(s;w))^2 \right]$$

def of $e_{ms}$

exchange $\nabla$ and $E$

$$= \frac{1}{2} E\left[ 2(V(s) - \hat{v}(s;w))' (+ \nabla \hat{v}(s;w)) \right]$$

true   approximate

$f(s)$ in linear case
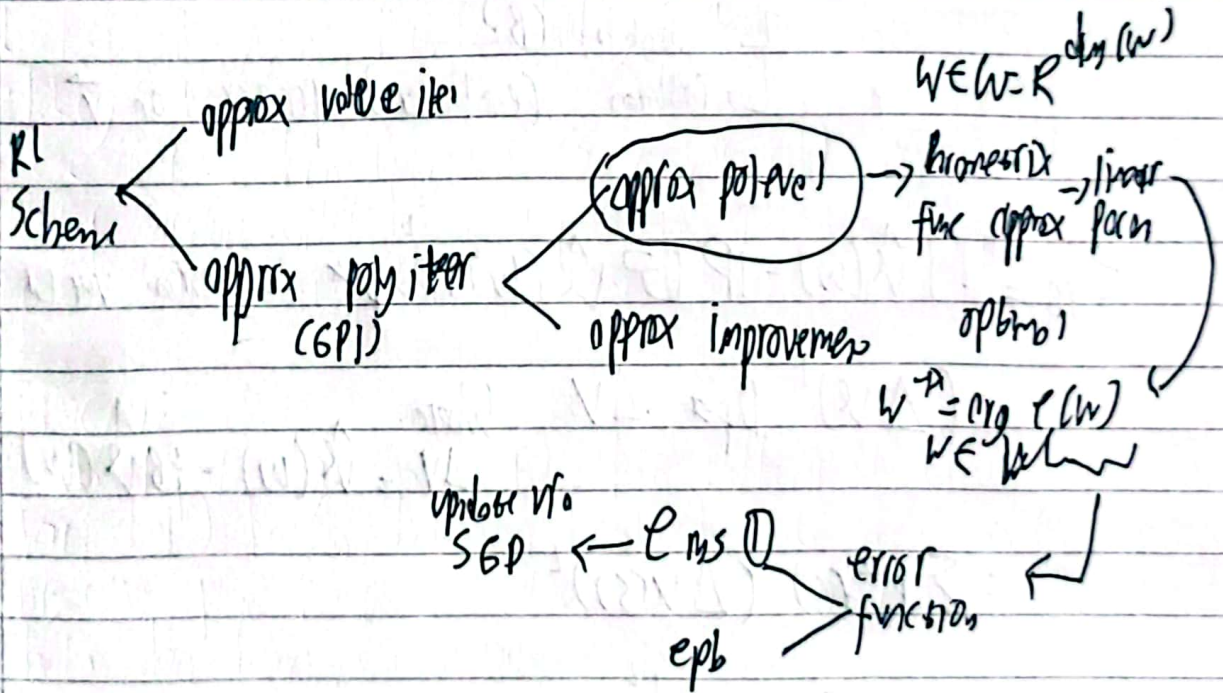
SL vs RL → Approximate the value (via TD, MC)

① training data    $V(s)$ is unknown (no training data)
exist

② iid data points    Markovian data (not iid) collected by the
(samples)    agent

$$= E\left[ [r(s,a) - g + \hat{v}(s;w)] - \hat{v}(s';w) ) f(s) \right]$$

approx $V(s)$ based BEE

$$\approx \left[ r(s,a) - \hat{g} + \hat{V}(s';w) - \hat{V}(s;w) \right] f(s)$$

↳ Sampling approx Using a Sample of $(s, a, r, 's')$
   to the expectation

approx value iter

$w \in w - R^{dm(w)}$

**RL Scheme**

approx poly iter (GPI)

(approx pol eval) → hrone matrix → linear
                                  func approx form

approx improvement        optimal

$$w^{*} = \arg \{ (w) $$
$$w \in \text{W} $$

update Vfo
SGD ← $\ell ms$ ① error function
                    epb

$\ell_{PB}$ ① doesn't include true value → $\Rightarrow$ Stable
   PB : Projected Bellman error
                ↳ Bellman operator

From BEE :          Immediate
Scalar :  $V(s) = r(s,a) - g + E[V(s')]$ ~ BEE
with
equality
deterministic                    Value next State )
Policy

Value(curr) = Func [ Value (next State )
                                   ↳T)

TD = RHS − LHS = 0 ... theory

Matrices: temporal
(vector)

$$\vec{\mathcal{V}} = \vec{r} - g\underline{1} + P \cdot \vec{\mathcal{V}} \iff (1-P)\mathcal{V} = r - g\underline{1}$$

$\overbrace{\qquad}^{\text{singlar}}$

↳ State-transition one step
State-transition $|S|$ by $|S|$

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

↳ $\mathrm{Matr} h\, bb(B)$
↳ Bellman (evaluator) operator on $\vec{V}$

$$\ell_{PB} \stackrel{\text{def}}{=} ||\hat{\mathcal{V}}(w) - P\, B\, \hat{\mathcal{V}}(w)||^2 \qquad P^* \text{ --- norm vector}$$

$$= (\Delta \mathcal{V})^T D_{P^*} \cdot \Delta V \qquad \text{where } \Delta V = \hat{\mathcal{V}}(w) - PB\hat{\mathcal{V}}(w)$$

$$= \sum P^*(s)\,(\Delta V(s))^2$$

$$= \mathbb{E}_{P^*}[\Delta V(s)]^2$$

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

$\underbrace{\qquad}_{\text{2-dim param}}$
2-dim param
vector

$B[\hat{\mathcal{V}}(w)]$ is not in Parameter space

↳ $w_2 \in R$ param $w$ cannot
represent $||B[\hat{V}\, w]$

$P\, ||B\, \hat{\mathcal{V}}\, ...$ inside param space

↳ to bring back $\hat{\mathcal{V}}(w)$
after $B[\hat{\mathcal{V}}(w)]$

The projection op subspaces $\mathbb{P}v = \boxed{F}\tilde{w}$ ← feature matrix

need derive P!

Where $\tilde{w} = \underset{w \in W}{argmin}\left\{\underbrace{\|Fw - v\|^2}_{\hat{v}} \cdot \underset{\downarrow}{p^\pi}\right\}$

$$\boxed{P = F(F^T D_{p\pi} F)^{-1} F^T D_{p\pi}}$$

determine $w^* = \underset{w}{argmin}\ \mathcal{E}(\theta)\ \ell_{PB}(w)$

How to get $w^*$

① take the grad with respect to $w$

$$\nabla \ell_{PB}(w)$$

② set the grad to zero

Then we can show that
$$w^* = X^{-1}Y$$

Where
$$X = \sum p^\pi(s) \sum p(s'|s) \left[ f(s)\{f(s) - f(s')\}^T \right]$$

$$= F^T D_{p\pi}(1-P)F$$

feature matrix      trans matrix

$$y(\theta) = \sum P^*(s) [r(s) \to \gamma 1] f(s)]$$
$$= F^T D_{P^*} (r - \gamma 1)$$

$$\hat{X} = \underbrace{\frac{1}{n} \sum_{i=1}^{n} F(s_i) [f(s_i) - F(s_i')]^T}_{\text{Sampling approx of } X}$$

Where $n$ is the number of Sample

$$\hat{y} = \underbrace{\frac{1}{n} \sum_{i=1} (r(s_i) - \gamma) f(s_i)}_{\text{Sampling approx of } Y}$$

$$\boxed{\hat{W} = \hat{X}^{-1} \hat{y}} \longrightarrow LSTD$$

$$\hat{x} \leftarrow \hat{x} + \underline{ET}$$

In RL, no separation between training vs testing

before convergence    after convergence

If he want to distinguish training & testing

in the Same environment

More challenge ⊖→ has Stationary env
⊖ differes but Similar env