- Stochastic gradient descent :-

$$\min_x f(x). \qquad \nabla f(x).$$

RGD: $\qquad x_{n+1} = x_n - \alpha_n \nabla f(x_n).$

(*) $\qquad \underset{\theta}{\min} \; \underset{x}{E}[f(x,\theta)].$

$$\downarrow$$

$$\underset{\theta}{\min} \; \boxed{\sum_x \mu(x) f(x,\theta).}$$

$$\boxed{*} \quad \theta_{n+1} = \theta_n - \alpha_n \left[ \sum_x \mu(x) \nabla_\theta f(x,\theta) \right]$$

$$\underbrace{\phantom{xxxxxxxxxxxxxxxxxxxxxxxx}}$$

$\boxed{*}$ Assuming $\mu(x)$ is known.

### Stochastic GD :-

$\boxed{*}$ At time instant 'n': you are given a sample

$$x_n \sim \mu(\cdot). \leftarrow \text{condition}$$

$$\boxed{\theta_{n+1} \leftarrow \theta_n - \alpha_n \left( \nabla_\theta f(x_n, \theta_n) \right).}$$

### Approximation :-

$$\boxed{Q(s,a)} :-$$

→ States are very high, infinite!!

$\qquad \qquad$ ↳ Mountain car :- $\boxed{\begin{array}{c} -0.6 \sim 2 \\ \end{array}}$

MDP
$$\downarrow$$
If model info is known $\Big\langle \begin{array}{l} \text{V.I} \\ \\ \text{P.I} \end{array}$

$$\downarrow$$

Model info is not known $\Big\langle \begin{array}{l} \text{Q-lear} \\ \text{SARSA} \\ \text{Exp sarsa} \end{array}$
(C.P.T.M Not known)

$$\downarrow$$

$\hookrightarrow$ + No. of states are $\Big\}$ → Approximation
$\qquad$ too high /
$\qquad$ infinite

**Prediction :-** given a policy $\pi$, want $V^{\pi}(\cdot)$

Parametrization : $V^{\pi}(s) \approx \hat{v}(s, \omega)$ $\rightarrow$ $\boxed{\omega \; \phi(s)}$

$\boxed{\text{coordinate}}$ $\uparrow$

$\boxed{\text{Parameter}}$ features corresponding to state $(s)$

feature $(s)$ = $\begin{bmatrix} d_1 \\ d_2 \\ d_3 \\ d_4 \\ \text{speed} \\ \text{acceleration} \end{bmatrix}$

$\boxed{\text{feature set}}$

- **Generalization :-**

  "If two states should have similar value fns, then their features should be closer".

  $V(S_1) = 10$     $V(S_2) = 10.2$

  Apprx: $\boxed{\omega^T \phi(S_1)}$     $\boxed{\omega^T \phi(S_2)}$

- "Approximated value function should be closer to the actual value function"

  $$J(\omega) = E\left[\left(\tilde{V}(s) - \hat{v}(s;\omega)\right)^2\right]$$

  $\boxed{S \sim M(\cdot)}$ $\rightarrow$ stationary distrn.

  $\pi$ : $\boxed{MDP} + \boxed{\pi} \rightarrow \boxed{MC}$ $\rightarrow$ stationary d$^n$. $\rightarrow \mu$

- $\omega_{n+1} = \omega_n + \frac{1}{2}\alpha\left[V(s) - \hat{v}(s,\omega)\right]\left(\nabla \hat{v}(s,\omega)\right). \leftarrow ?$

approximate $V(s)$ and use an estimate of $V(s)$ in above eqn.

Monte-carlo estimate:-

$$V^{\pi}(s) = \mathbb{E}_{\tau}\left[G_{\tau}\right].$$

gradient M.C algorithm:-

$$\omega_{n+1} = \omega_n + \frac{\alpha}{\alpha}\left[G_{\tau} - \hat{v}(s,\omega)\right]\nabla\hat{v}(s,\omega). \leftarrow$$

At time $n$:- generate a trajectory $\tau$:

compute $G_{\tau}$:-

$$\boxed{S_1, a, 1, S_2, a_2, 2, S_3, a_3, 4, S_4.} \leftarrow$$

$\boxed{G_{\tau}}$: $1 + \gamma(2) + \gamma^2(4)$

total discounted sum.

TD(0) estimate :-

$$\boxed{V^{\pi}(s)} = \mathbb{E}\left[\underline{r} + \gamma\,\underline{V^{\pi}(s')}\right] \quad \text{unbiased} \\ \longrightarrow \text{estimate of } V^{\pi}(\cdot)$$

$$\boxed{s' \sim P(\cdot\,|\,s, a\pi)}.$$

semi gradient TD(0) update rule:-

$$\checkmark \boxed{\omega_{n+1}} \leftarrow \omega_n + \alpha\left[r + \gamma\,V^{\pi}(s') - \hat{v}(s,\omega)\right]\nabla\hat{v}(s,\omega).$$

$$\longrightarrow \hat{v}(s',\omega).$$

$$\boxed{\text{Not known !}}$$

$$\boxed{\omega_{n+1} \leftarrow \omega_n + \alpha\left[r + \gamma\,\hat{v}(s',\omega) - \hat{v}(s,\omega)\right]\nabla\hat{v}(s,\omega)} \leftarrow$$

$$\boxed{n : \boxed{(\underline{s}, \underline{a}, \underline{r}, \underline{s'})}.}$$

$$\boxed{\hat{v}(s,\omega) = \omega^T\phi(s)}$$

**Policy gradient :-**

(control)

ee construct an ==Objective== function which can ==be optimized== to obtain ==optimal policy== ①

1. **Parameterization of policy:**

$$\pi(S) \approx \pi(S, \theta)$$

→ Parameter to be optimized.

$$\pi_\theta(S, a) = \dfrac{e^{\theta^T \phi(S,a)}}{\sum_b e^{\theta^T \phi(S, b)}}$$

one example

$\theta$ $\phi(S,a)$ → features corr. to state (S) & action (a)

---

Bellman Eqn
↓ Model free
SARSA

Q-Bellman
↓ Model free
→ Q-learning.

---

$\pi(S, a_1) \rightarrow$
$\pi(S, a_2) \rightarrow$
$\vdots$
$a_n$

$$a^* = \arg\max_a \; Q(S, a)$$

---

2. **constructing the objective function :-**

$(\theta)$

$$J(\theta) = V_{\pi_\theta}(S_0).$$

→ Find $\theta$ s.t Value fn. Starting from State $S_0$ is maximum.

$S_0 \leftarrow$ fixed initial state

$V(S) \leftarrow$ value fn. corresponding to State S.

---

3. can be optimized:

$\nabla J(\theta)$    ML.   TD(0).

$$\nabla J(\theta) \propto \sum_S \mu(S) \sum_a q_\pi(S,a) \cdot \nabla \pi(S, a)$$

$$E_{S \sim \mu(S)} \left[ \sum_a q_\pi(S,a) \cdot \nabla \pi(S,a) \right].$$

$$\theta^* \leftarrow \arg\max_\theta J(\theta)$$

$$\theta_{n+1} \leftarrow \theta_n + \alpha_n \boxed{\nabla J(\theta_n)}$$

2y needs to be estimated

$$\boxed{S \sim \mu(\cdot).}$$

1. computing : $\nabla \pi(S, \theta).$

2. Approximating : $q_\pi(S, a)$

$$\begin{cases} MC \longrightarrow \text{Reinforce} \\ TD(0) \longrightarrow \text{Actor-critic} \end{cases}$$

3. sampling lechin : $S \sim \mu(\cdot)$
   $\hookrightarrow$ not known.

general rule :
$$\theta_{t+1} = \theta_t + \alpha \left( \sum_a \hat{q}(S,a) \cdot \nabla \pi(S,a) \right).$$

$$\boxed{\textcircled{*} \quad \nabla \log \pi(S,a) = \boxed{\frac{\nabla \pi(S,a).}{\pi(S,a).}}}$$

$$\sum_S \mu(S) \sum_a q_\pi(S,a) \cdot \nabla \pi(S,a).$$

$$= \sum_S \mu(S) \sum_a \cdot q_\pi(S,a) \cdot \frac{\pi(S,a).}{\pi(S,a)} \cdot \nabla \pi(S,a).$$

$$= \sum_S \mu(S) \sum_a \pi(S,a) \left[ q_\pi(S,a) \cdot \nabla \log \pi(S,a) \right].$$

$$= \mathbb{E}_{(S,a)} \left[ q_\pi(S,a) \cdot \nabla \log \pi(S,a) \right]$$

$$\boxed{\theta_{t+1} \leftarrow \theta_t + \alpha \left[ \hat{q}(S,a) \cdot \nabla \log \pi(S,a) \right].}$$

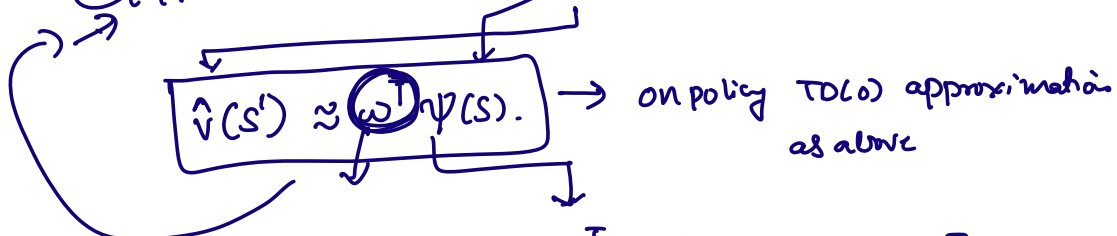$\hookrightarrow$ final update rule.

Reinforce :-

$$\theta_{t+1} \leftarrow \theta_t + \alpha \left[ \boxed{q_t} \nabla \log \pi(S,a) \right].$$

TD(0) :- approximation

   ↳ Actor-critic :-

$$Q(s,a) = \mathbb{E}\left[r + \gamma v^*(s')\right].$$

$$\theta_{t+1} \leftarrow \theta_t + \alpha \left[ (r + \gamma v^*(s'))\ \nabla \log \pi(s,a)\right].$$

$$\hat{v}(s') \approx \omega^T \psi(s). \quad \rightarrow \text{ on policy TD(0) approximation}$$
$$\text{as above}$$

$$\theta_{t+1} \leftarrow \theta_t + \alpha \left[ (r + \gamma \omega^T \psi(s))\ \nabla \log \pi(s,a)\right].$$

         ↳ final AC algorithm.

2 updates :   $\omega$ : estimating v.f $\rightarrow$ critic $\Big\}$
        $\theta$ : improving policy $\rightarrow$ Actor.