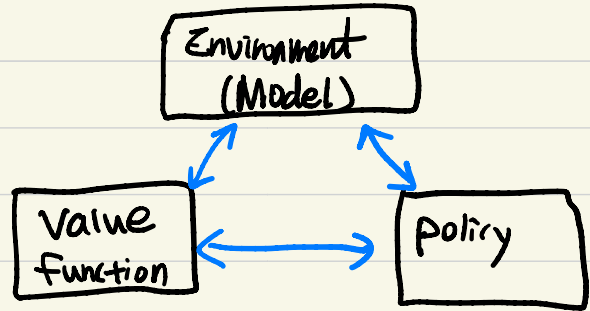# Lec 07 a — Policy Optimization.

**\*where we are.**

① Markov Decision Process

② Dynamic Programming
   - Learn value function
   - Implicit policy

③ Policy Optimization
   - No value function
   - Learn policy
     ex) Policy gradient



---

Formulation :

$$\max_{\theta} \ \mathbb{E}\left[ \sum_{t=0}^{T-1} \gamma^t \, r(s_t, a_t) \,\Big|\, \pi_\theta \right] \qquad - (1)$$
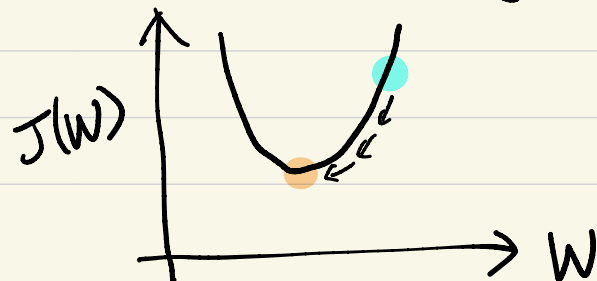
(1) is finite-horizon discounted problem

The goal is to maximize the return, $R_t$ in an episodic setting.
Parameterized policy $\pi_\theta(a|s)$ is used for stochastic policy.

Q: How to solve this optimization problem?
    A: Gradient-based technique.

Gradient Descent Algorithm.



$$\text{"} \ W_{k+1} = W_k - \alpha \nabla_w J(w) \ \text{"}$$

update law.

- gradient of objective fcn.

$$g = \nabla_\theta \mathbb{E}\left[ \sum_{t=0}^{T-1} \gamma^t r(S_t, a_t) \mid \pi_\theta \right]$$

$$\theta_{K+1} = \theta_K + \alpha g.$$

- example, Direct Policy grad.

$$\mathbb{E}\left[ \sum_{t=0}^{1} r(S_t, a_t) \mid \pi_\theta \right] = \mathbb{E}\left[ \underbrace{r(S_0, a_0)}_{①} + \underbrace{r(S_1, a_1)}_{②} \mid \pi_\theta \right]$$

①: $\nabla_\theta \mathbb{E}\left[ r(S_0, a_0) \right] = \nabla_\theta \int \underbrace{r(S_0, a_0)}_{RV} \underbrace{\mu(S_0) \pi_\theta(a_0 \mid S_0)}_{prob.} \underline{\underline{dS_0}}$

②: $\nabla_\theta \mathbb{E}\left[ r(S_1, a_1) \right] = \nabla_\theta \int r(S_1, a_1) \pi_\theta(a_1 \mid S_1) P(S_1 \mid S_0, a_0)$

$$\mu(S_0) \pi_\theta(a_0 \mid S_0) \underline{\underline{dS_0}} \; \underline{\underline{dS_1}}$$

Every subsequent term adds additional dimension of integration.
⟺ It's computationally intractable to compute gradient
analytically

We are going to "estimate" the gradient "g".

⟺ policy gradient.