# Policy Gradient Theorem:

problem: order a product more or less every day.

reward depends on $P$: $\Delta$ element
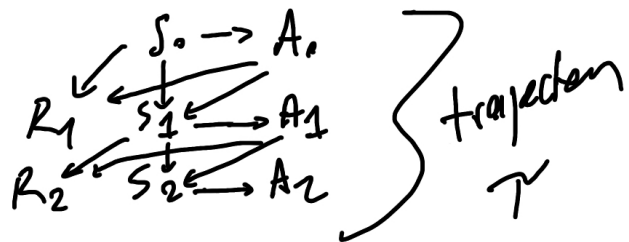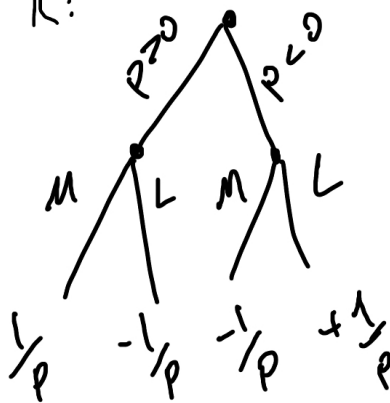
$a = \{$ order more, order less $\}$

$s$: $\{$ change in demand $\}$ $(P)$

(or $P$)

( if $P \to +$ more demand

if $P \to -$ less demand )

$R$:



$$P(\text{more}) = \sigma(\theta_p)$$

→ sigmoid

→ learn this one

→ demand change



Mapping state of the world $(P)$ to the probability of taking a certain action.

↳ Policy $(\pi)$

find ideal $\theta$ in

$\frac{\partial J}{\partial \theta} = \frac{\partial}{\partial \alpha} E_\theta(R)$

P(More) : $\delta(\theta_P)$

$J(\theta) = E(\theta)$ [total rewd which we want to maximize over all $\gamma$]

$\hookrightarrow \frac{\partial}{\partial \theta} \sum_T \pi_\theta(T) R(T)$

$\frac{\partial J}{\partial \theta} \rightarrow$ gradient ascent

$\parallel$

$\sum_T \pi_\theta(T) \frac{\partial}{\partial \theta} \left[ \log \pi_\theta(\gamma) \right] R(\gamma)$

log bc $\frac{\partial}{\partial \theta} \pi_\theta(\gamma)$ is long.

$\frac{\partial}{\partial \theta} \log \pi_\theta(\gamma) = \frac{\partial}{\partial \theta} \underbrace{\log P(S_0)}_{0} +$

$\sum_{t=0}^{T} \frac{\partial}{\partial \theta} \log \underbrace{P(A_t | S_t)}_{\text{policy basically}}$

no need to model the environment. model-free

$+ \sum_{t=0}^{\gamma} \frac{\partial}{\partial \theta} \underbrace{\log P(R_{t+1}, S_{t+1} | S_t, A_t)}_{0}$

$$E_\theta \left( R(\tau) \cdot \sum_{t=0}^{T} \frac{\partial}{\partial \theta} \log P(A_t \mid S_t) \right) \Rightarrow$$

$$\frac{\partial J}{\partial \theta} = p \, P(\text{loss})$$