

Actor - Critic

Hyperparameters:

- Initial policy params θ
- Policy representation (ANN? Linear?)
- Actor step size α
- Critic step size β .
- Value function representation, v_w .
- Initial value function weights w .

$$v^\pi(s) = \mathbb{E} \left[\sum_{k=0}^{\infty} \gamma^k R_{t+k} \mid S_t = s; \pi \right]$$

For each episode:

For each time t

Agent observes S_t

Agent selects action A_t using π_θ

Env responds with S_{t+1} and R_t

$$\delta_t = R_t + \gamma v_w(S_{t+1}) - v_w(S_t) \quad // \text{TD-error}$$

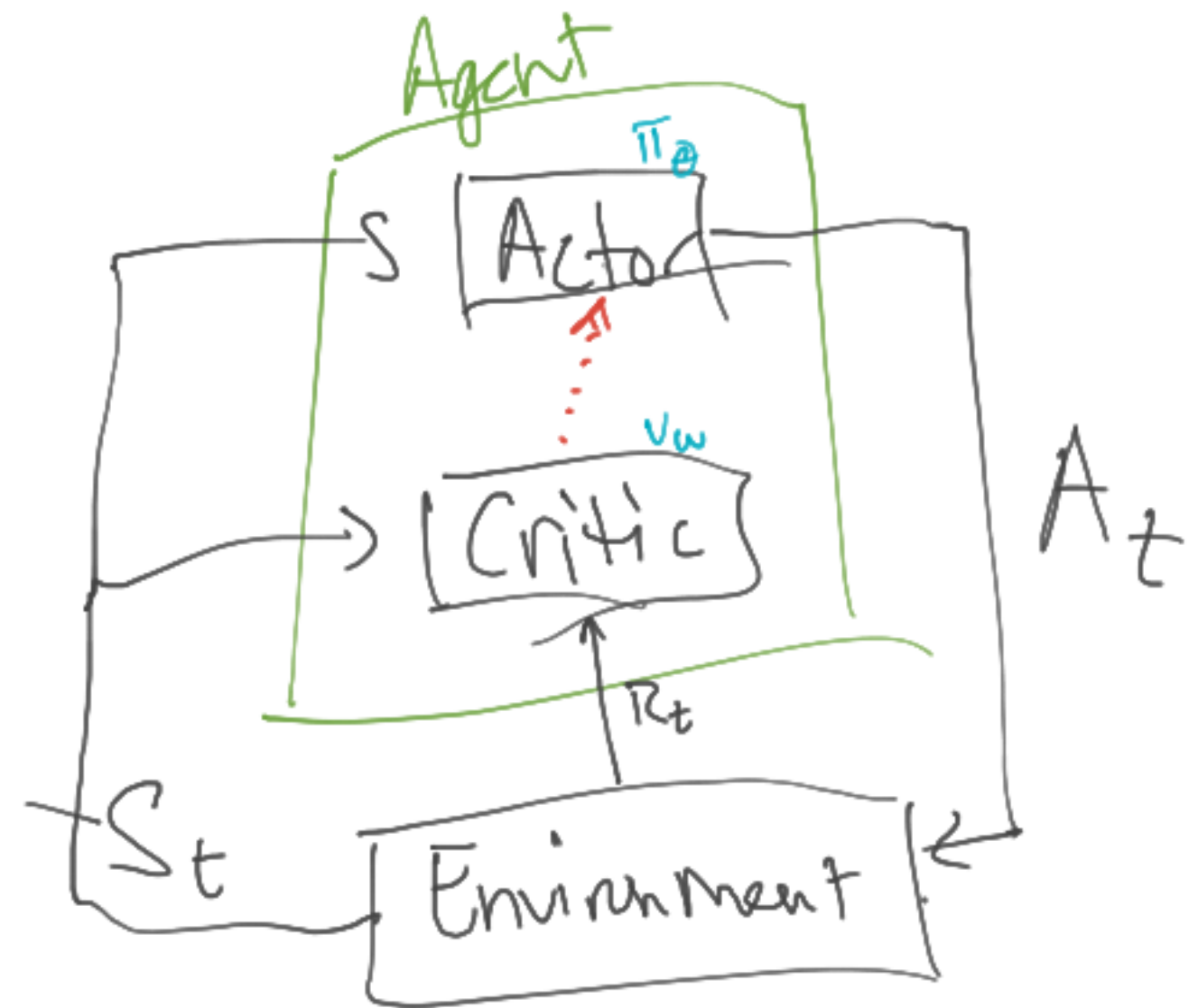
$$\theta_i \leftarrow \theta_i + \alpha \delta_t \frac{\partial \ln(\pi_\theta(S_t, A_t))}{\partial \theta_i} \quad // \text{Actor update}$$

$$w_j \leftarrow w_j + \beta \delta_t \frac{\partial v_w(S_t)}{\partial w_j} \quad // \text{critic update.}$$

For each time t

~~$$w_j \leftarrow w_j + \beta \left(\sum_{k=0}^{\infty} \gamma^k R_{t+k} - v_w(S_t) \right) \frac{\partial v_w(S_t)}{\partial w_j} \quad // \text{critic update}$$~~

δ_t



Actor - Critic

Hyperparameters:

- Initial policy params θ
- Policy representation (ANN? Linear?)
- Actor step size α
- Critic step size β .
- Value function representation, v_w .
- Initial value function weights w .

For each episode:

For each time t

Agent observes S_t

Agent selects action A_t using π_θ

Env responds with S_{t+1} and R_t

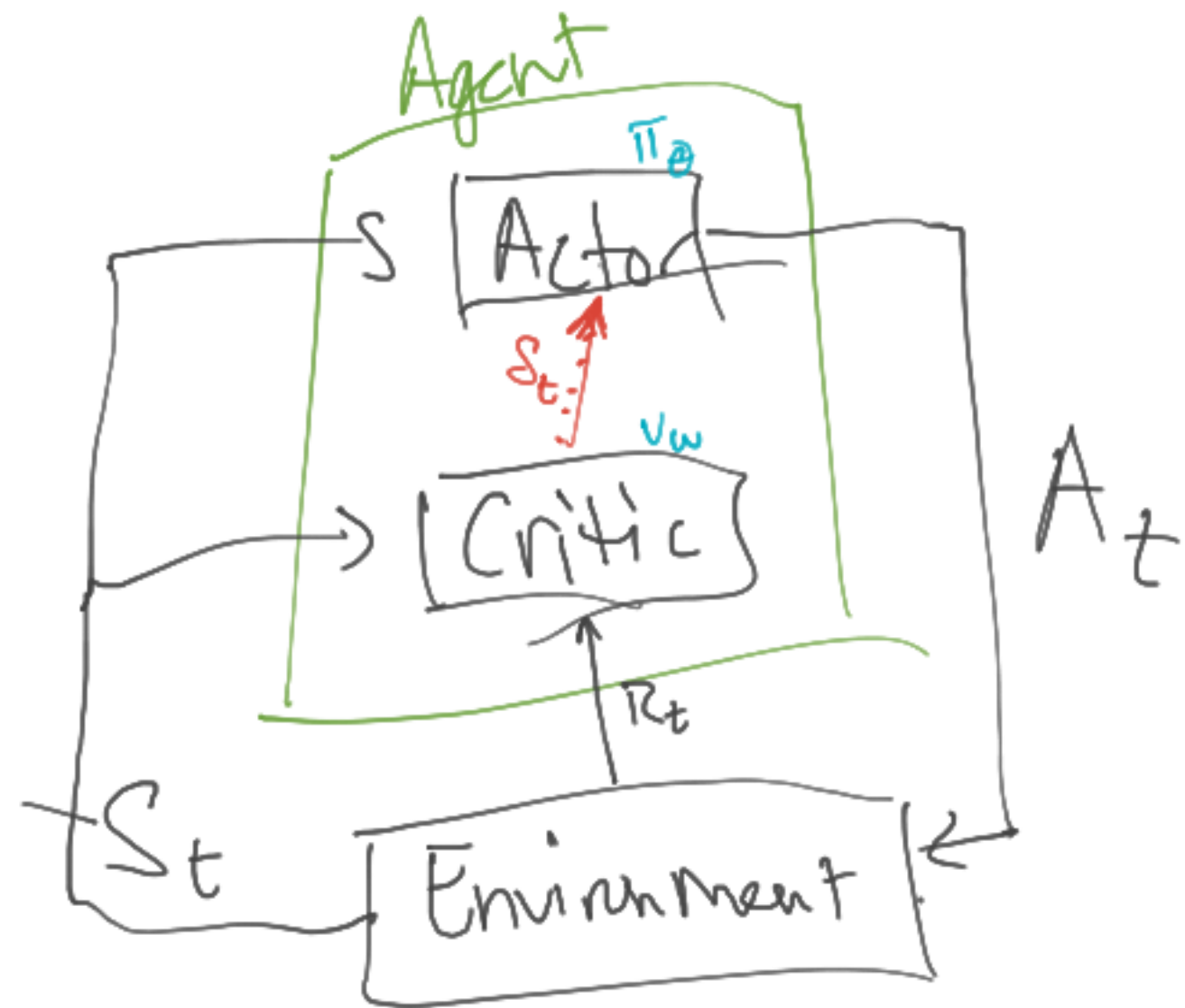
$$\delta_t = R_t + \gamma v_w(S_{t+1}) - v_w(S_t) \quad // \text{TD-error}$$

$$v_i, \theta_i \leftarrow \theta_i + \alpha \gamma \delta_t \frac{\partial \ln(\pi_\theta(S_t, A_t))}{\partial \theta_i} \quad // \text{Actor update}$$

$$v_j, w_j \leftarrow w_j + \beta \delta_t \frac{\partial v_w(S_t)}{\partial w_j} \quad // \text{critic update.}$$

Theory says to include
In practice it is bad. Almost nobody includes this term.

$$v^\pi(s) = \mathbb{E} \left[\sum_{k=0}^{\infty} \gamma^k R_{t+k} \mid S_t = s; \pi \right]$$



Learning v^π

- Supervised learning problem!
→ regression

input	output
s_0	$\sum_{k=0}^{\infty} \gamma^k R_{0+k}$
s_1	$\sum_{k=0}^{\infty} \gamma^k R_{1+k}$
\vdots	

Parametric model (f_w), v_w . $v_w(s)$ is an estimate of $v^\pi(s)$
Recall: (grad descent on least squares loss) (pointwise like backprop)

$w_{k,j}$ in old notation is w_j here.

$$w_j \leftarrow w_j + \alpha \left(\sum_{k=0}^{\infty} \gamma^k R_{t+k} - \cancel{f_w(x)} \right) \frac{\partial \cancel{f_w(x)}}{\partial w_j}$$

$v_w(s_t)$

$$w_j \leftarrow w_j + \alpha \left(\sum_{k=0}^{\infty} \gamma^k R_{t+k} - v_w(s_t) \right) \frac{\partial v_w(s_t)}{\partial w_j}$$

linear func approx

$$v_w(s_t) = \sum_j w_j \phi_j(s_t)$$
$$w_j \leftarrow w_j + \alpha \left(\sum_{k=0}^{\infty} \gamma^k R_{t+k} - v_w(s_t) \right) \phi_j(s_t)$$

Moving critic update into the episode:

$$V_j \quad w_j \leftarrow w_j + \alpha \left(\sum_{k=0}^{\infty} \gamma^k R_{t+k} - v_w(S_t) \right) \frac{\partial v_w(S_t)}{\partial w_j}$$

$$V_j \quad w_j \leftarrow w_j + \alpha \delta_t \frac{\partial v_w(S_t)}{\partial w_j}$$

$$R_t + \gamma R_{t+1} + \gamma^2 R_{t+2} \dots$$

This is a (stochastic) gradient update on least squares loss.

....., S_t, A_t, R_t, S_{t+1}

NOT gradient update. (Not grad descent or least squares) (can diverge!)

$$R_t + \gamma (R_{t+1} + \gamma R_{t+2} + \dots)$$

δ_t $v_w(S_{t+1})$ prediction

$$V_j \quad w_j \leftarrow w_j + \alpha \left(\underbrace{R_t}_{\text{label}} + \gamma \underbrace{v_w(S_{t+1})}_{\text{prediction}} - v_w(S_t) \right) \frac{\partial v_w(S_t)}{\partial w_j}$$

Policy Gradient

Run gradient ascent on $J(\theta) = \mathbb{E} \left[\sum_{t=0}^{\infty} \gamma^t R_t; \pi_{\theta} \right]$

$$\forall j, \theta_j \leftarrow \theta_j + \alpha \frac{\partial J(\theta)}{\partial \theta_j}$$

Algs presented are approximately doing this. Hence, called "policy gradient".

→ REINFORCE

→ Standard Actor-Critic.

→ Many AC algorithms.

→ This is just one.

→ wait until episode end to update, use $\sum_{k=0}^{\infty} \gamma^k R_{t+k}$ not θ

→ Missing γ^t term
→ Include in Do updates simultaneously

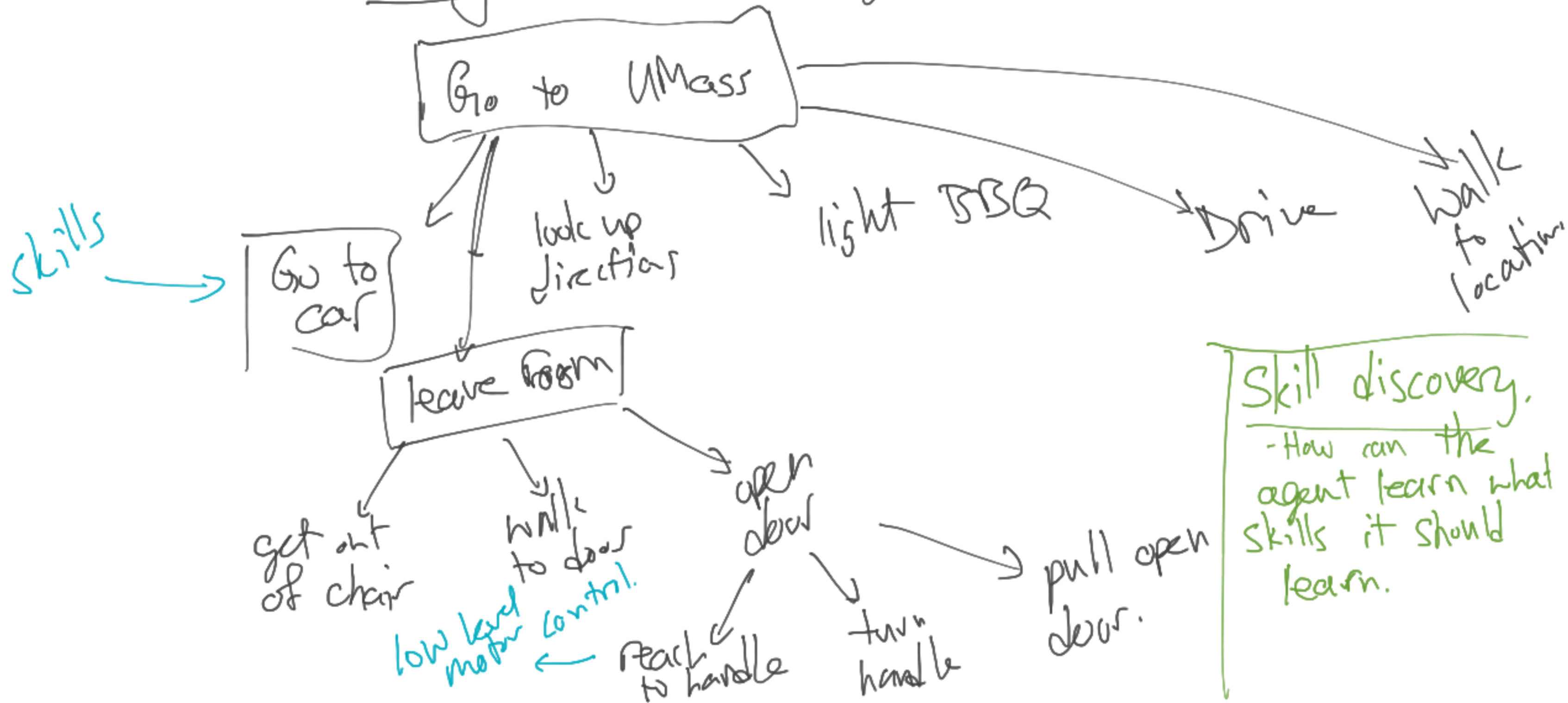
Not an exams or HW.

Advanced Topics in RL

"Hierarchical RL"

Skills or options.

A Skill is a policy that the agent can choose to run.



Skill discovery.
- How can the agent learn what skills it should learn.

off-policy evaluation

- Have data from current policy
- Goal: predict (with confidence bound)
how good a new policy would be if
we used it instead, but without actually
using it.