

ECE-GY 6143

Intro to ML

D Recap: Unsupervised learning

D Reinforcement learning (RL)

O Model-free RL

O Policy gradients.

Supervised learning

- Regression
- Classification

Unsupervised learning

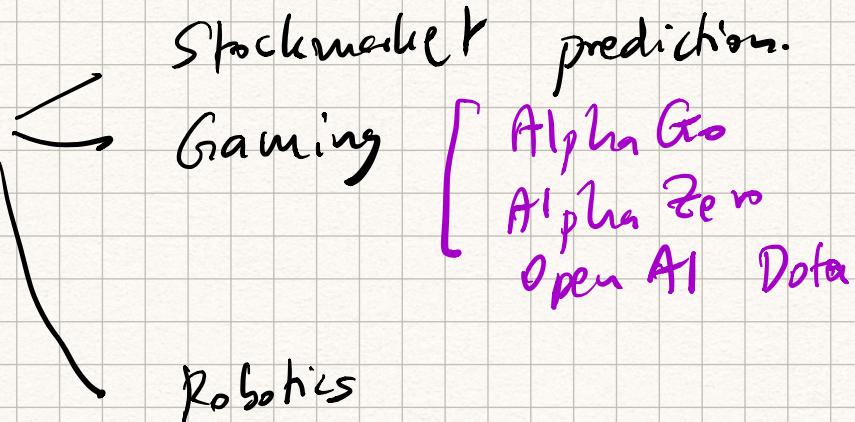
- Data visualization
- Dimensionality reduction.

- Dataset is static
- Learning is performed offline
- Predictions do not affect data samples

Reinforcement learning (RL)

- Dataset is dynamically changing.
- Learning is online
- Explicit modeling of feedback loops.

## Applications



[ RL  $\leftrightarrow$  Control theory  
Dynamical systems, optimal nonlinear control ] -

---

## Setting

- Environment
  - Agent / state
  - Actions at time step
  - Actions  $\rightarrow$  rewards
  - Goal: Maximize rewards.
- 

$$s_{t+1} = f(s_t, a_t)$$

$\hookrightarrow$  state transition function / dynamical system .

$$\text{Reward} = r(s_t, a_t)$$

Goal:  
Agent observes

$\tilde{s}_t = (s_0, a_0, s_1, a_1, \dots, \underline{s_t})$

↳ Trajectory / rollout ↳ Horizon

$r(s_t, a_t)$  for  $t = 0 \dots (T-1)$

Predict a strategy / policy for deciding the next action.

Policy:  
Function  $\pi$   $\rightarrow \underline{a_t} = \underline{\pi(x_t)}$

Apply ML to solve this prediction problem!

→ Choose representation for  $\pi$

→ Choose a loss function

$$R. = \sum_{t=0}^{T-1} -\underline{r}(s_t, a_t)$$

subject to

$$\underline{s_{t+1}} = f(\underline{s_t}, \underline{a_t})$$

$$\underline{a_t} = \underline{\pi(x_t)}, \quad s_0 \text{ is given.}$$

→ Optimize this function.

Linear models  
logistic  
neural nets  
Deep RL

Need to figure out

$$\boxed{\frac{\partial R}{\partial \pi}}$$

??

- Idea: Assume that  $\pi$  is not deterministic, instead it predicts a prob. distribution of actions.

$a_t$  random

$s_t$  random

i.  $R$  is a random variable.

Modified problem

min

$$E_{\pi(c)} \underbrace{R(c)}_{= E - \sum_{t=0}^{T-1} r(a_t, s_t)}$$

$$\text{s.t. } s_{t+1} = f(s_t, a_t) \\ a_t = \pi(s_t)$$

Training samples

→ Multiple trajectories / rollouts.

Supervised

$\rightarrow R_L$

Unsupervised.

$x \rightarrow y$

only 1  $y$  per rollout( $x$ )

$x$

## Policy gradients.

- Q-learning
- Value iterations
- Actor-Critic.

Assume  $\pi$  is a linear policy parameterized by  $\theta$ .

$$\underbrace{\frac{\partial}{\partial \theta} E_{\pi(\tau)} R(\tau)} = ?$$

Log-derivative trick.

$$\frac{\partial \log \pi(\tau)}{\partial \theta} \Rightarrow \frac{1}{\pi(\tau)} \cdot \frac{\partial}{\partial \theta} \pi(\tau).$$

$$\therefore \underbrace{\frac{\partial}{\partial \theta} \pi(\tau)} = \pi(\tau) \underbrace{\frac{\partial}{\partial \theta} \log \pi(\tau)}$$

$$\underbrace{\frac{\partial}{\partial \theta} E_{\pi(\tau)} R(\tau)} \Rightarrow \frac{\partial}{\partial \theta} \sum_{\tau} \pi(\tau) R(\tau)$$

$$= \sum_{\tau} R(\tau) \frac{\partial}{\partial \theta} \pi(\tau)$$

$$= \sum_{\tau} \underbrace{\pi(\tau)}_{\text{R}(\tau)} \underbrace{R(\tau) \frac{\partial}{\partial \theta} \log \pi(\tau)}_8$$

$$\nabla E_{\pi(\tau)} R(\tau) = E_{\pi(\tau)} \left[ \underbrace{R(\tau)}_{\text{R}(\tau)} \frac{\partial}{\partial \theta} \log \pi(\tau) \right]$$

Gradient of expected rewards

= Expected value of  $\underbrace{\text{Reward} \times \log \text{derivative}}$ .

∴ Just like in SGD,

- sample different rollouts
- Compute approximation to the gradient of expected rewards.
- Update  $\theta$  in that direction.

Repeat

1) Sample rollout  $\tau = (\overbrace{s_0, a_0, s_1, \dots, s_T})^T$

2) Compute  $R(\tau) = \sum_{t=0}^{T-1} -r(s_t, a_t)$

3)  $\underline{\theta} \leftarrow \underline{\theta} - \gamma \frac{\partial R(\tau)}{\partial \underline{\theta}} \cancel{\log \pi_\theta(\tau)}$

REINFORCE

↳  $R$  appears but never its gradient.

↳ Useful when rewards are discrete.

↳ Environment  $f$  does not appear,

only the rollout  $\pi_{\text{random}}(\tau)$ .

"Model-free reinforcement learning".

↳ "Model-based RL"



Derivative-free optimization

Gradient descent:

$$\theta \leftarrow \theta - \eta \nabla R(\theta)$$

Random search:

Derivative-free optimization.

- 1) Pick a random vector  $v$
- 2) Compute step size ( $\eta$  either negative or positive) that minimizes  $R(\theta + \eta v)$ .
- 3)  $\theta \leftarrow \theta + \eta v$ .