

Reinforcement Learning and Control

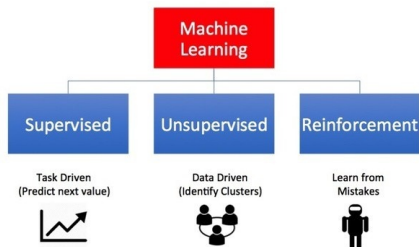
Luis Pimentel

ECE 4803 Mathematical Foundations of Data Science

November 24, 2020

Reinforcement Learning: Machine Learning Perspective

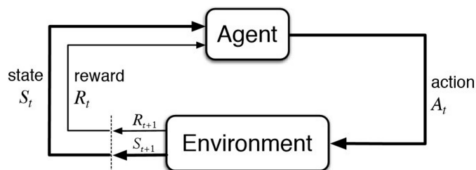
Types of Machine Learning



Goal of Reinforcement Learning

Learn to how behave in an unknown and uncertain enviroment using **data** from **previous behavior**.

Definitions



- **Agent**: system that performs actions affecting its state
- **Environment**: the world in which the agent performs
- **State S_t** : particular condition of an agent at a specific time
- **Action A_t** : activity performed by the agent
- **Policy π_t** : state to action mapping (**what you are trying to learn!**)
- **Reward R_t** : feedback from the environment as a result of action

Optimization Problem Formulation

- Let τ represent a state-action trajectory s.t.
 $\tau = [\mathbf{s}_0, \mathbf{a}_0, \mathbf{s}_1, \mathbf{a}_1, \dots; \mathbf{s}_{N-1}, \mathbf{a}_{N-1}, \mathbf{s}_N]$
- Let $f(\mathbf{s}_t, \mathbf{a}_t, e_t)$ represent the dynamics of the system with Gaussian noise e_t .

$$\underset{\pi}{\text{maximize}} \quad \mathbb{E}[R(\tau)] \qquad R(\tau) = \sum_{t=0}^N r(\mathbf{s}_t, \mathbf{a}_t, t)$$

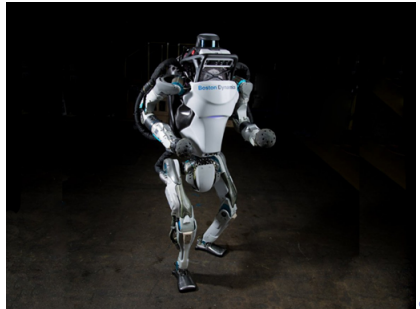
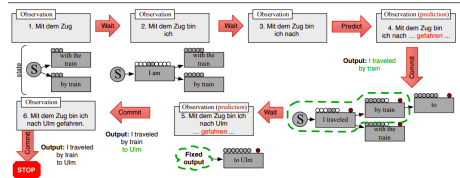
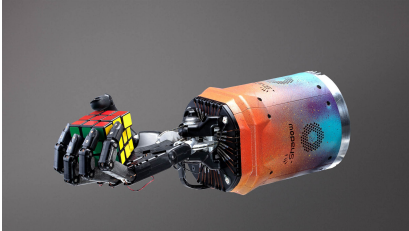
subject to

$$\mathbf{s}_{t+1} = f(\mathbf{s}_t, \mathbf{a}_t, e_t)$$

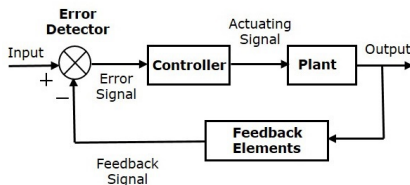
$$\mathbf{a}_t = \pi(\tau)$$

- Our goal is to optimize our policy such that we maximize our reward.
- We set up our reward function such that
 - “good” action leads to a “good” state, thus a positive reward
 - “bad” action, which leads to a “bad” state, thus a negative reward

Reinforcement learning applications:



Control Systems

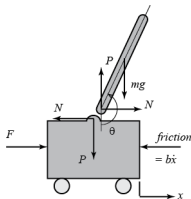


$$\dot{\mathbf{x}}(t) = \mathbf{A}\mathbf{x}(t) + \mathbf{B}\mathbf{u}(t)$$

$$\mathbf{y}(t) = \mathbf{C}\mathbf{x}(t)$$

- $\dot{\mathbf{x}}$: differential equation describing how the state changes
- \mathbf{x} : state of the system
- \mathbf{u} : input to the system
- \mathbf{y} : output of the system
- $\mathbf{A}, \mathbf{B}, \mathbf{C}$: state and control transition matrix, state observation matrix

Example: Inverted Cart-Pendulum Linearized Model



$$\begin{aligned}(I + ml^2)\ddot{\phi} - mgl\phi &= ml\ddot{x} \\ (M + m)\ddot{x} + b\dot{x} - ml\ddot{\phi} &= u\end{aligned}$$

$$\begin{bmatrix} \dot{x} \\ \ddot{x} \\ \dot{\phi} \\ \ddot{\phi} \end{bmatrix} = \begin{bmatrix} 0 & 1 & 0 & 0 \\ 0 & \frac{-(I+ml^2)b}{I(M+m)+Mml^2} & \frac{m^2gl^2}{I(M+m)+Mml^2} & 0 \\ 0 & 0 & 0 & 1 \\ 0 & \frac{-mlb}{I(M+m)+Mml^2} & \frac{mgl(M+m)}{I(M+m)+Mml^2} & 0 \end{bmatrix} \begin{bmatrix} x \\ \dot{x} \\ \phi \\ \dot{\phi} \end{bmatrix} + \begin{bmatrix} 0 \\ \frac{I+ml^2}{I(M+m)+Mml^2} \\ 0 \\ \frac{ml}{I(M+m)+Mml^2} \end{bmatrix} u$$

What happens when we don't have a model??

Can we still control a system without understanding the detailed mathematics and physics behind it?

We can reformulate the control problem as a Reinforcement Learning problem and learn the system input policy.

- $\tau = [\mathbf{x}_0, \mathbf{u}_0, \mathbf{x}_1, \mathbf{u}_1, \dots; \mathbf{x}_{N-1}, \mathbf{u}_{N-1}, \mathbf{x}_N]$
- We parameterize our policy by θ .

$$\underset{\theta}{\text{maximize}} \quad \mathbb{E}[R(\tau)]$$

$$R(\tau) = \sum_{t=0}^N r(\mathbf{x}_t, \mathbf{u}_t, t)$$

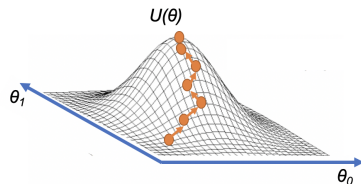
subject to

$$\mathbf{x}_{t+1} = \dot{f}(\mathbf{x}_t, \mathbf{u}_t, e_t)$$

$$\mathbf{u}_t = \pi(\tau, \theta)$$

Solving the optimization problem

We can solve this problem using our favorite method: **gradient ascent**!



$$\underset{\theta}{\text{maximize}} \ J(\theta)$$

Algorithm 1: Gradient Ascent for Policy Parameterization

Result: θ

initialize: $\theta = \theta_0$, γ : learning rate

while *unconverged* **do**

$\nabla_{\theta} J(\theta_i)$ = compute_policy_gradient(θ_i)
 $\theta_{i+1} = \theta_i + \gamma \nabla_{\theta} J(\theta_i)$

- If we don't have a model of our system, how do we calculate the policy gradient?
 - Solution: we can use input and output data to create an approximation of the gradient.
 - Apply an input, observe the output and reward
- Estimating the Policy Gradient is a major research area in Reinforcement Learning.
- Two methods of Policy Gradient estimation:
 - Finite-Difference Methods
 - Episodic REINFORCE

Finite Difference Method

Main idea: Apply a perturbation $\delta\theta$ to our θ parameter over many rollout trajectories to create an estimation of the Policy Gradient.

1)

$$J(\theta + \delta\theta_m) = J(\theta) + \nabla_{\theta} J(\theta)^T \delta\theta_m$$

$$J(\theta + \delta\theta_m) - J(\theta) = \nabla_{\theta} J(\theta)^T \delta\theta_m$$

$$\Delta\hat{J}_m = \nabla_{\theta} J(\theta)^T \delta\theta_m$$

2) set up a Least-Squares Regression Problem:

$$\begin{bmatrix} \Delta\hat{J}_1 \\ \Delta\hat{J}_2 \\ \vdots \\ \Delta\hat{J}_m \end{bmatrix} = \begin{bmatrix} \delta\theta_1^T \\ \delta\theta_2^T \\ \vdots \\ \delta\theta_m^T \end{bmatrix} \nabla_{\theta} J(\theta)$$

$$\Delta\hat{\mathbf{J}} = \Delta\mathbf{\Theta} \nabla_{\theta} J(\theta)$$

$$\nabla_{\theta} J(\theta) = \left(\Delta\mathbf{\Theta}^T \Delta\mathbf{\Theta} \right)^{-1} \Delta\mathbf{\Theta}^T \Delta\hat{\mathbf{J}}$$

Finite Difference Method Algorithm

$$\nabla_{\theta} J(\theta) = \left(\Delta \Theta^T \Delta \Theta \right)^{-1} \Delta \Theta^T \Delta \hat{J}$$

Algorithm 2: Finite Difference Gradient Estimator

Result: θ

initialize: $\theta = \theta_0$, γ : learning rate, m : rollouts

σ^2 : variance, ϵ : Gaussian noise

while *unconverged* **do**

for m rollouts **do**

$\delta \theta_m = \sigma \epsilon$

$\Delta \hat{J}_m = J(\theta + \delta \theta_m) - J(\theta)$

$\nabla_{\theta} J(\theta_i) = \left(\Delta \Theta^T \Delta \Theta \right)^{-1} \Delta \Theta^T \Delta \hat{J}$

$\theta_{i+1} = \theta_i + \gamma \nabla_{\theta} J(\theta_i)$

Finite Difference Method Considerations

- Very efficient for deterministic systems.
- Different configurations:
 - How to generate the perturbation $\delta\theta$
 - How to compute $\Delta\hat{J}$
 - How to choose convergence criterion
 - Tuning of other hyper-parameters
- Have been successfully implemented on robotics systems.
- Require close supervision of robotics engineer as badly generated $\delta\theta$ can lead to de-stabilization and failure in optimization.

Episodic REINFORCE

- In Episodic REINFORCE we consider path probabilities as the basis of our objective function.
- $\tau = (\mathbf{x}_0, \mathbf{u}_0, \dots, \mathbf{x}_{N-1}, \mathbf{u}_{N-1}, \mathbf{x}_N)$
- $R(\tau)$ is the accumulated cost over a trajectory
- $p(\tau)$ represents the path probability of the trajectory, which using Bayesian and Markov properties can be expressed as:

$$p(\tau) = \underbrace{p(\mathbf{x}_0)}_{\text{init. dist.}} \prod_{i=0}^N \underbrace{p(\mathbf{x}_{i+1} | \mathbf{x}_i, \mathbf{u}_i)}_{\text{state trans. dist.}} \underbrace{p(\mathbf{u}_i | \mathbf{x}_i; \theta)}_{\text{policy } \pi \text{ dist.}}$$

$$R(\tau) = \sum_{t=0}^N r(\mathbf{x}_t, \mathbf{u}_t, t)$$

Optimization Formulation

$$\text{maximize}_{\theta} J(\theta) = \text{maximize}_{\theta} \int p(\tau) R(\tau) d\tau$$

The expected return of the policy can be written as the sum of the expected rewards over all trajectories.

$$\nabla_{\theta} J(\theta) = \int \nabla_{\theta} p(\tau) R(\tau) d\tau$$

$$\nabla_{\theta} \log(p(\tau)) = \frac{1}{p(\tau)} \nabla_{\theta} p(\tau)$$

$$\int p(\tau) \nabla_{\theta} \log(p(\tau)) R(\tau) d\tau$$

$$\nabla_{\theta} J(\theta) = E_{p(\tau)} \left[\nabla_{\theta} \log(p(\tau)) R(\tau) \right]$$

$$\nabla_{\theta} J(\theta) = E_{p(\tau)} \left[\nabla_{\theta} \log \left(\underbrace{p(\mathbf{x}_0)}_{\text{init. dist.}} \prod_{i=0}^N \underbrace{p(\mathbf{x}_{i+1} | \mathbf{x}_i, \mathbf{u}_i)}_{\text{state trans. dist.}} \underbrace{p(\mathbf{u}_i | \mathbf{x}_i; \theta)}_{\text{policy } \pi \text{ dist.}} \right) R(\tau) \right]$$

$$\nabla_{\theta} J(\theta) = E_{p(\tau)} \left[\sum_{i=0}^N \nabla_{\theta} \log \left(\underbrace{p(\mathbf{u}_i | \mathbf{x}_i; \theta)}_{\text{policy } \pi \text{ dist.}} \right) R(\tau) \right]$$

Gradient of the expected reward only depends the expected reward over a trajectory and the log probability of the parameterized control policy.

The derivatives with respect to the control system do not need to be computed to estimate the gradient.

REINFORCE Algorithm

$$\nabla_{\theta} J(\theta) = \frac{1}{M} \sum_{m=0}^M \left[\left(\sum_t \nabla_{\theta} \log(p(\mathbf{u}_{m,t} | \mathbf{x}_{m,t}; \theta)) \right) \left(\sum_t r(\mathbf{x}_{m,t}, \mathbf{u}_{m,t}, t) \right) \right]$$

Algorithm 3: REINFORCE Algorithm

Result: θ

initialize: $\theta = \theta_0$, γ : learning rate, m : rollouts

while *unconverged* **do**

for m rollouts **do**

 └ Sample trajectories by running the policy.

$\nabla_{\theta} J(\theta) = E_{\tau; \theta}$

$\theta_{i+1} = \theta_i + \gamma \nabla_{\theta} J(\theta_i)$

REINFORCE Considerations

- Guaranteed to converge to the true gradient at fastest theoretical rate.
- Policy parameter variations are no longer needed eliminating computationally expensive process and potentially endangering the policy gradient estimation.
- In real world implementations a single rollout is enough to get a good estimation of the policy gradient.

Application to Control: LQR Optimal Control

Optimal Control Problem Formulation:

$$\dot{\mathbf{x}} = \mathbf{A}\mathbf{x} + \mathbf{B}\mathbf{u}$$

$$\mathbf{u} = -\mathbf{K}\mathbf{x}$$

$$J(\mathbf{x}, \mathbf{u}) = \mathbf{x}^T \mathbf{Q} \mathbf{x} + \mathbf{u}^T \mathbf{R} \mathbf{u}$$

Control policy can be solved analytically with simple knowledge of the system dynamics and cost function:

$$\mathbf{K} = \mathbf{R}^{-1} \mathbf{B}^T \mathbf{P}$$

$$\mathbf{P}\mathbf{A} + \mathbf{A}^T \mathbf{P} + \mathbf{Q} - \mathbf{P}\mathbf{B}\mathbf{R}^{-1} \mathbf{B}^T \mathbf{P} = \dot{\mathbf{P}}$$

What happens when we don't have a dynamics model?

Can we use Reinforcement Learning to learn the optimal control gain K and apply the feedback controller?

We can reformulate as a Reinforcement Learning problem:

$$\underset{\theta}{\text{maximize}} \quad \mathbb{E}[R(\tau)]$$

$$R(\tau) = \sum_{t=0}^N r(\mathbf{x}_t, \mathbf{u}_t, t)$$

$$r(\mathbf{x}_t, \mathbf{u}_t, t) = -\mathbf{x}^T \mathbf{Q} \mathbf{x} - \mathbf{u}^T \mathbf{R} \mathbf{u}$$

$$\mathbf{u} = -\theta \mathbf{x}$$