

• CME 241 Notes : Reinforcement Learning for Predictions (17 - Feb - 2021)

- No more required readings beyond this chapter
  - Module 3 : Reinforcement Learning
- 
- DP/ADP assume access to a probability model (knowledge of  $P_{\theta}$ )
    - But if you know probabilities, you don't need to interact with environment
  - In RL, we don't have access to a probability function
    - ⇒ You must interact with actual environment
      - In state, make move, and environment produces one instance of next state and reward (individual Experiences)
      - Often you try to build a gated  $P_{\theta}$  function by supervised learning
  - Rather than trial and error, you may want to fit a model by supervised learning.
    - The real world is not stationary
    - Models must be constantly updated
    - Also, state space is huge (curse of dimensionality)
    - Curse of Modeling: Probability estimation on large state space is intractable
      - It's difficult to stitch together individual probability models
        - \* It's more practical to stitch samples together
          - This is much more reasonable. Gives sampling model for  $P_{\theta}$ , which serves as a proxy for the real world.
      - Let RL interact with simulated environment
        - Allows you to run over many experiences

- RL Needs an interface to sample experiences
  - ↳ RL learns the value function w/o access to probability, but rather by accessing a stream of simulated experiences
  - ↳ Under the right settings, RL will converge on value function
  - ↳ Like real life, RL does not attempt to build a probability model.
  - ↳ RL is a trial and error approach to linking actions to accumulated rewards (returns).
    - ↳ How do you disentangle actions from eventual returns?
    - ↳ Actions have highly overlapping reward sequences.
    - ↳ Some actions result in delayed reward.
- Our attention will shift to the Q-value function
  - ↳ Expected return conditional on state and action.
  - ↳ RL constantly updates Q-value function
    - ↳ Supposed to be an expectation. Updated through new samples
    - ↳ Updates done through function approximation
      - ↳ Linear Function
      - ↳ Neural Network
    - ↳ Almost all of RL is based on Bellman Equations
- Prediction is a stepping stone to control
  - ↳ Today and Friday we will talk about RL for prediction

- Prediction: Evaluate MDP with fixed policy  $\pi$  using individual experiences

$$\text{MDP} + \pi = \text{MRP}$$

↳ We access samples from MRP to estimate value function

$$s_0 \rightarrow r_1, s_1 \rightarrow r_2, s_2 \rightarrow \dots$$

"Atomic Experience"

"Trace Experience"

$$(s, r, s') \equiv \text{Transition Step.}$$

- Software Design: Interface can be an iterable of atomic experiences

↳ Or you could have a stream of trace experiences by sticking together transition steps.

↳ Iterable can be either a Sequence or an Iterator (Stream)

#### • Monte Carlo Prediction

• You can get many sample returns for a state and you can fit a model with Supervised learning.

↳ Implemented with incremental manner in FunctionApprox.Update() method.

↳ You can only do an update at the end of the trace experience because we are talking about returns.

↳ Returns are easier to calculate backwards

↳ Read "batch-batch" code

- ↳ Minimizes MSE of predicted and observed return
- ↳ "Episode" is a finite trace experience
- ↳ MC Method requires every trace experience to end.

- Code:

"trs" is a stream of episodes

Specify an initial estimate of the initial function approximation

- Prediction output produces estimate of updated function approximations using a generator
- Function approximation updates() calls update again and again. Sequence of updates for all episodes.
- Update formula:

$$\Delta w = \alpha (G_t - V(s_t; w)) \cdot \nabla_w V(s_t; w)$$

↑                      ↑                      ↑  
 Learning rate      Return residue      Gradient of Expected returns

- Tabular MC Prediction

- Denote  $V_n(s_i)$  as the estimate of VF for first  $n$  occurrences of  $s_i$

$$V_n(s_i) := V_{n-1}(s_i) + f(n) \cdot (Y_i^{(n)} - V_{n-1}(s_i))$$

↑  
Count to weight function

↳ One default approach is  $f(n) = \frac{1}{n}$

↳ With this approach, the update function is a running average

- ↳ For tabular MC, incremental averaging done separately for each state
- ↳ Tabular MC does incremental average

- ↳ Law of Large Numbers  $\Rightarrow$  Average approaches true expectation
- ↳ Tabular MC is a specialized case of linear function approx
  - ↳ Feature matrix is an  $n$ -dimensional identity matrix
  - ↳ Weight vector is average of each state.
  - $\therefore Y_i^{(n)} = W_i^{(n)}$  is the gradient of the loss function
- In general, you want to use a time-decaying weighted average of returns.
  - ↳ Lower weights for older data.
- This version of MC is called "Each Visit MC"
  - ↳ Update for each occurrence of a given state
  - ↳ Opposed to "First Visit MC" = update only for first visits
- MC produces unbiased estimates, but it can be slow to converge
- MC requires complete trace experiences
  - ↳ The inability to use incomplete experiences is a key disadvantage
- Temporal Difference Prediction: Begin with a tabular view
  - ↳ TD Exploits recursive structure of VF in Bellman Equation
    - ↳ Replace  $G_t$  with  $R_{t+1} + \gamma V(S_{t+1})$  using atomic experiences
    - ↳ This is a bootstrapping approach
  - ↳ Update to VF happens after every experience
  - $R_{t+1} + \gamma V(S_{t+1})$  = "Temporal Difference Target"
  - $\delta = R_{t+1} + \gamma V(S_{t+1}) - V(S_t)$  = "Temporal Difference Error"

$$V(s_t) \leftarrow V(s_t) + \alpha r$$

- ↳ Key idea: Unlike MC, TD is non-episodic
- ↳ We may have continuous traces  $\Rightarrow T \leq 1$
- ↳ TD can be run on any stream of atomic experiences
  - ↳ You can split up / rearrange atomic experiences
- Switching to TD Prediction with Function Approximation