

• 27 Jan 2021 Notes: Dynamic Programming Continued.

• Bellman Optimality Operator :

$$B^{Q(v)}(v) = B^*(v)$$

↑
Action prescribed by
Greedy policy is
the one which satisfies
Bellman Optimality Equation

B^* motivated by MDP Bellman Policy Equation

Bellman Optimality Operator Equation can be expressed as

$$V^* = B^*(V^*)$$

∴ V^* is a fixed point

B^* is a contraction function

→ Read Proof from slide 19 and textbook,
will likely be similar to exam.

• Steps of proof :

① Monotonicity Property

$$x \geq y \Rightarrow B^*(x) \geq B^*(y)$$

② Constant Shift Property

$$B^*(x + c) = B^*(x) + yc$$

$$\Rightarrow \max_{s \in N} |(B^*(x) - B^*(y))(s)| \leq \dots \text{ see slide 19.}$$

- Value iteration algorithm

- Complexity is $O(n^2 k)$
- Repeatedly apply B^* on π until convergence

- Optimal Policy from Optimal Value Function

$$V^{\pi^*} = V^*$$

Searching for policy which achieves optimal value

\hookrightarrow Use greedy policy function G

$$B^{G(v^*)}(v^*) = B^*(v^*)$$

\Downarrow

v^*

$$\therefore B^{G(v^*)}(v^*) = v^*$$

$\hookrightarrow v^*$ is a fixed point

\hookrightarrow We know at $v^{G(v^*)}$ $B^{G(v^*)}$ has a unique fixed point

$$\therefore V^{G(v^*)} = v^*$$

Evaluating MDP with greedy policy at v^* achieves v^*

$$\therefore \pi^* = G(v^*)$$

\uparrow

We know how
to get this

- Slide 21 has a good demonstration of algorithm.
 - ↳ Rows are policy evaluations
 - ↳ Repeatedly apply Bellman Operator
 - ↳ Between each row is Policy Improvement Step.
- Two notions of consistency:
 - ① Across rows, make $\text{VF } V$ consistent with / close to V^* of the policy π
 - ② Between rows, make π consistent with / close to Greedy policy $G(V)$ of $\text{VF } V$

↳ Achieving one type of consistency loses the other type of consistency. Therefore, alternate between them
- Generalized Policy Iteration is the foundation of RL
 - ↳ Intuition: lie partially between the two notions of consistency
- Value iteration schematic on slide 25
 - ↳ "Lazy" version of policy iterations (type on slide)
- Large scale RL uses Partial policy evaluation and policy improvement
 - # Policy iteration is more precise, but slower than value iteration
 - ↳ Both work for a small state space
 - ↳ Value iterations can have convergence issues
 - ↳ But for finite state spaces, both work and value iteration is better.

- To implement DP algorithms, see slide 26
 - ↳ Asynchronous DP algorithm uses prioritized sweeping, ranked based on gaps between states in the algorithm.
 - ↳ Address large gaps first.
 - ↳ You need links from destination states to start states to sort queue (disadvantage).
- Batteries Break... Continue from slide 27 (~40 minute point of lecture)
- Finite Horizon MDP:
 - ↳ Each sequence terminates in a finite number of time steps T
 - ↳ Each time step has a separate (from other time steps) set of states
 - ↳ Walking back in time
 - ↳ MDP with fixed policy is equivalent to MRF
 - ↳ Backward Induction Algorithm: $O(m^2 T)$
 Decrement t from T to 0 and calculate V_t^* from V_{T+1}^*
 - ↳ Backward induction applies to both prediction and control
 - ↳ For control, extract optimal policy using argmax
- Lots of problems in finance show up as finite horizon problem
- E.g. Dynamic Pricing: End-of-life / End-of-season clearance
 - ↳ How much to reduce price at end of season/life
 - T days away from end of season

M units of inventory

Maximize expected total sales revenue over T days

- ↳ Choose discrete prices P_t
- ↳ Each price associated with Poisson Distribution of demand
 - ↳ Price elasticity of demand

Visualize as grid of time and inventory

- ↳ See slide 32 for model

→ How do you model λ in Poisson distribution of demand based on pricing.

• Generalizations to non-tabular Algorithm

- ↳ Value function represented as a table in tabular algorithms
- ↳ You cannot do this for infinite states
- ↳ Instead, collect a sample of states and apply Bellman operator for those states. Create a function approximation with those samples (using interpolation).
- ↳ Fundamental principles are the same.
- ↳ Known as Approximate Dynamic programming

• Approximate Dynamic Programming (New slide deck)

• Function Approximation:

- ↳ Predictor variable $x \in \mathcal{X}$, response $y \in \mathbb{R}$

$$P(y|x) = f(x; \omega)(y)$$

↑
parameters

→ Assume data in form $[(x_i, y_i) \mid 1 \leq i \leq n]$ is empirical distribution D

→ We want to construct model probability distribution M

→ Requires D and M to minimize cross-entropy loss:

$$H(D, M) = - \mathbb{E}_D [l_M]$$

→ Minimize using gradient descent

- First method in FunctionApprox is to solve for ω given pairs of x, y
- Linear Function Approximation : slides 7-9
- Deep Neural Network : slide 10