# ECE 5984
# Model-Free Reinforcement Learning
# Temporal Difference Methods

Jason J. Xuan, Ph.D.

Department of Electrical & Computer Engineering
Virginia Tech

# Temporal-Difference Learning

- TD methods learn directly from episodes of experience

- TD is *model-free*: no knowledge of MDP transitions / rewards

- TD learns from *incomplete* episodes, by *bootstrapping*

- TD updates a guess towards a guess

# MC and TD

- Goal: learn $v_\pi$ online from experience under policy $\pi$
- Incremental every-visit Monte-Carlo
    - Update value $V(S_t)$ toward *actual* return $G_t$

$$V(S_t) \leftarrow V(S_t) + \alpha \left( G_t - V(S_t) \right)$$

- Simplest temporal-difference learning algorithm: TD(0)
    - Update value $V(S_t)$ toward *estimated* return $R_{t+1} + \gamma V(S_{t+1})$

$$V(S_t) \leftarrow V(S_t) + \alpha \left( R_{t+1} + \gamma V(S_{t+1}) - V(S_t) \right)$$

    - $R_{t+1} + \gamma V(S_{t+1})$ is called the *TD target*
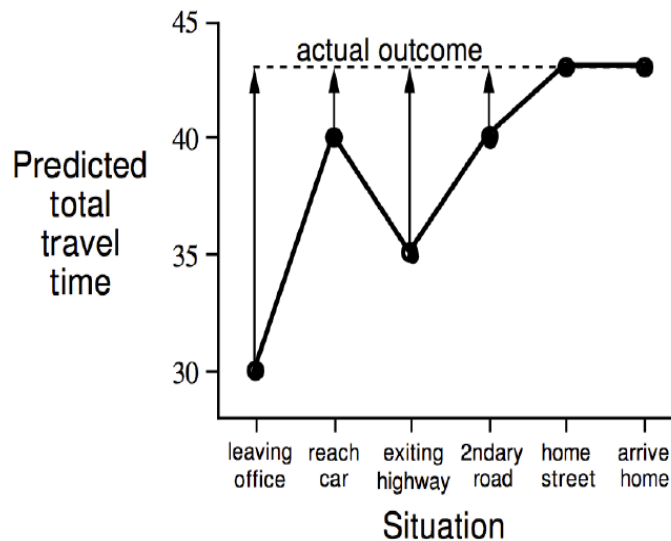    - $\delta_t = R_{t+1} + \gamma V(S_{t+1}) - V(S_t)$ is called the *TD error*

# Example: Driving Home

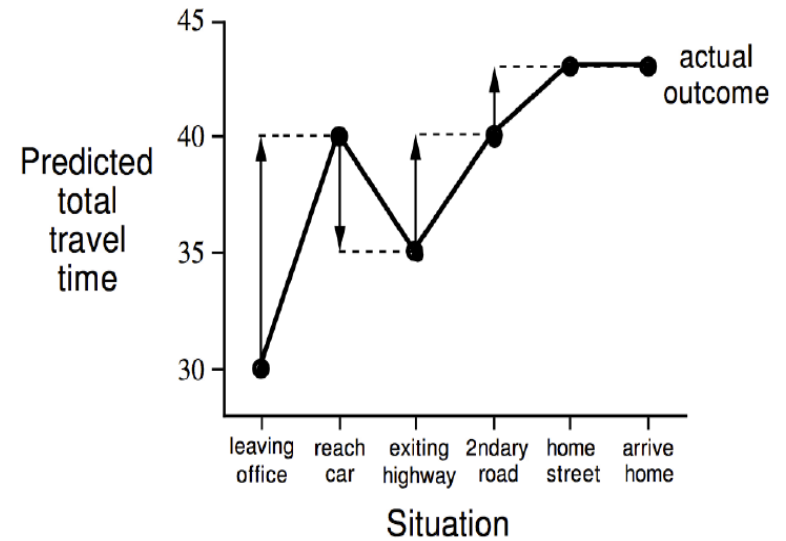| State | Elapsed Time (minutes) | Predicted Time to Go | Predicted Total Time |
|---|---|---|---|
| leaving office | 0 | 30 | 30 |
| reach car, raining | 5 | 35 | 40 |
| exit highway | 20 | 15 | 35 |
| behind truck | 30 | 10 | 40 |
| home street | 40 | 3 | 43 |
| arrive home | 43 | 0 | 43 |

# Driving Home: MC vs TD



Changes recommended by Monte Carlo methods ($\alpha=1$)

Changes recommended by TD methods ($\alpha=1$)

# Pros and Cons: MC vs TD

- TD can learn *before* knowing the final outcome
  - TD can learn online after every step
  - MC must wait until end of episode before return is known
- TD can learn *without* the final outcome
  - TD can learn from incomplete sequences
  - MC can only learn from complete sequences
  - TD works in continuing (non-terminating) environments
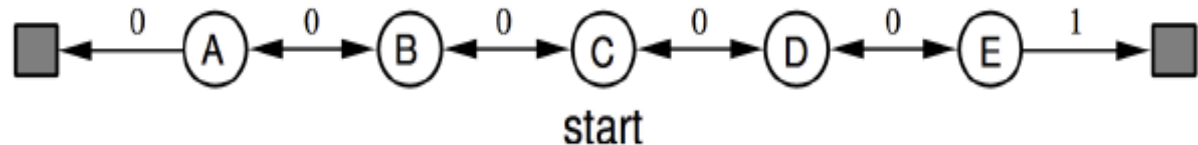  - MC only works for episodic (terminating) environments

# Bias-Variance Tradeoff

- Return $G_t = R_{t+1} + \gamma R_{t+2} + ... + \gamma^{T-1} R_T$ is *unbiased* estimate of $v_\pi(S_t)$

- True TD target $R_{t+1} + \gamma v_\pi(S_{t+1})$ is *unbiased* estimate of $v_\pi(S_t)$

- TD target $R_{t+1} + \gamma V(S_{t+1})$ is *biased* estimate of $v_\pi(S_t)$

- TD target is much lower variance than the return:
    - Return depends on *many* random actions, transitions, rewards
    - TD target depends on *one* random action, transition, reward
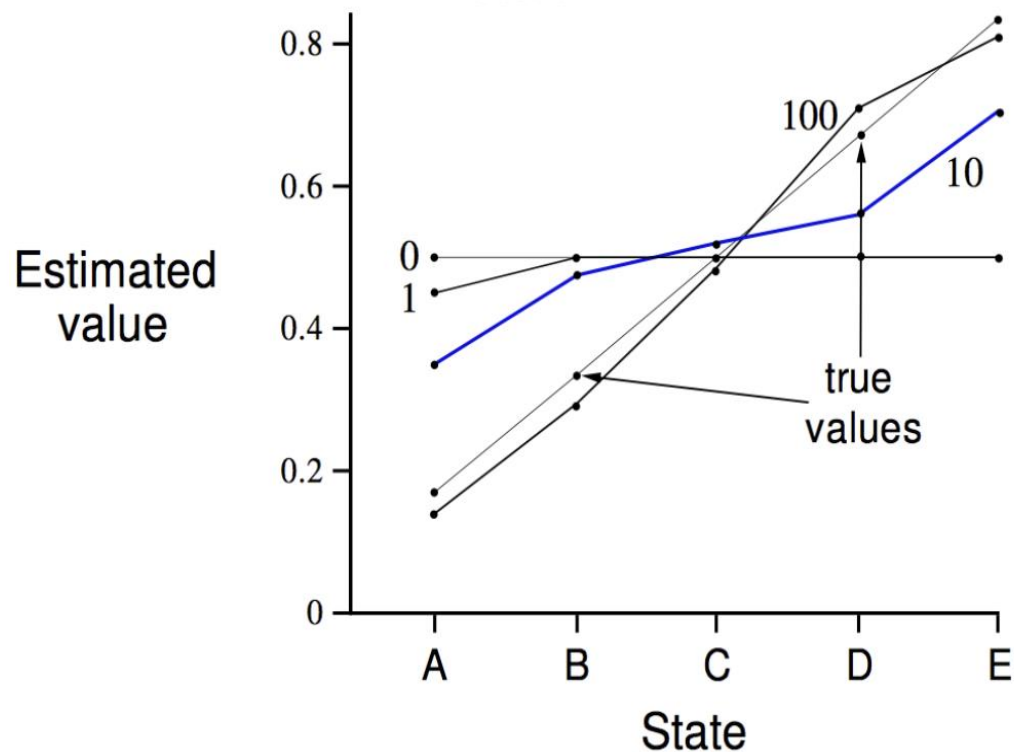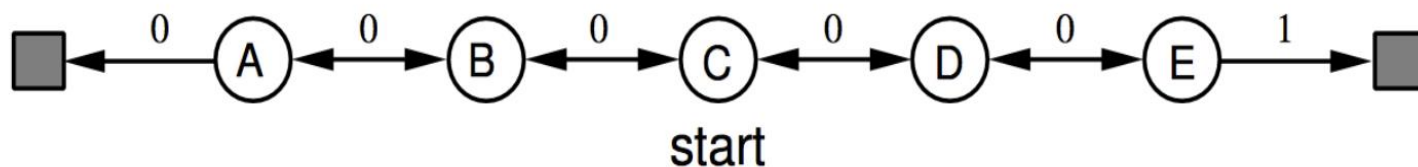
- MC has high variance, zero bias
    - Good convergence properties

    - Not very sensitive to initial value
    - Very simple to understand and use
- TD has low variance, some bias
    - Usually more efficient than MC
    - TD(0) converges to $v_\pi(s)$

    - More sensitive to initial value
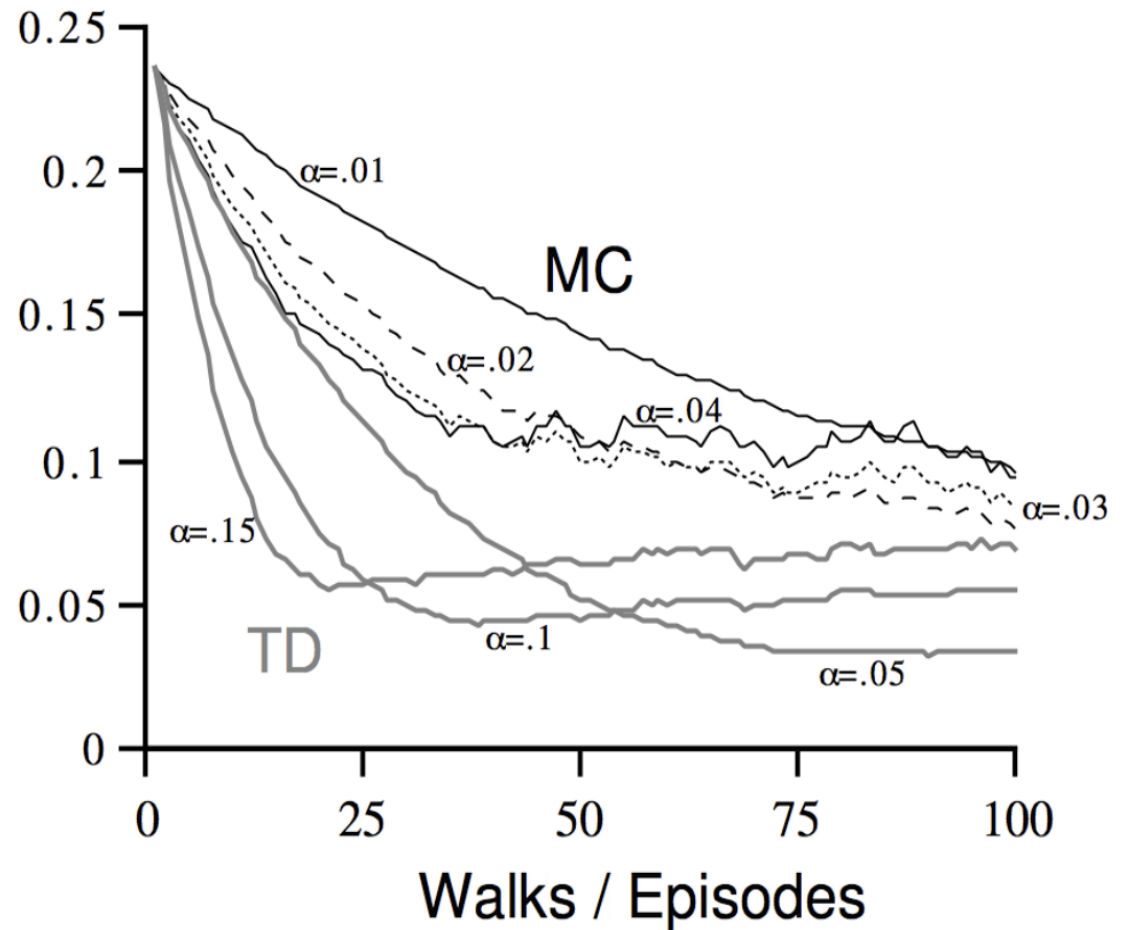
# Example: Random Walk



- All episodes start in the center state C
- proceed either left or right by one state on each step, with equal probability
- Episodes terminate either on the extreme left or the extreme right.
- When an episode terminates on the right a reward of 1 occurs; all other rewards are zero.
- Because this task is undiscounted and episodic, the true value of each state is the probability of terminating on the right if starting from that state.
- The true values of all the states, A through E, are 1/6, 2/6, 3/6, 4/6, 5/6

# Random Walk: MC vs TD

# Certainty Equivalence

- MC converges to solution with minimum mean-squared error
  - Best fit to the observed returns

$$\sum_{k=1}^{K}\sum_{t=1}^{T_k}\left(G_t^k - V(s_t^k)\right)^2$$

- TD(0) converges to solution of max likelihood Markov model
  - Solution to the MDP $\langle \mathcal{S}, \mathcal{A}, \hat{\mathcal{P}}, \hat{\mathcal{R}}, \gamma \rangle$ that best fits the data

$$\hat{\mathcal{P}}_{s,s'}^a = \frac{1}{N(s,a)}\sum_{k=1}^{K}\sum_{t=1}^{T_k}\mathbf{1}(s_t^k, a_t^k, s_{t+1}^k = s, a, s')$$

$$\hat{\mathcal{R}}_s^a = \frac{1}{N(s,a)}\sum_{k=1}^{K}\sum_{t=1}^{T_k}\mathbf{1}(s_t^k, a_t^k = s, a)r_t^k$$
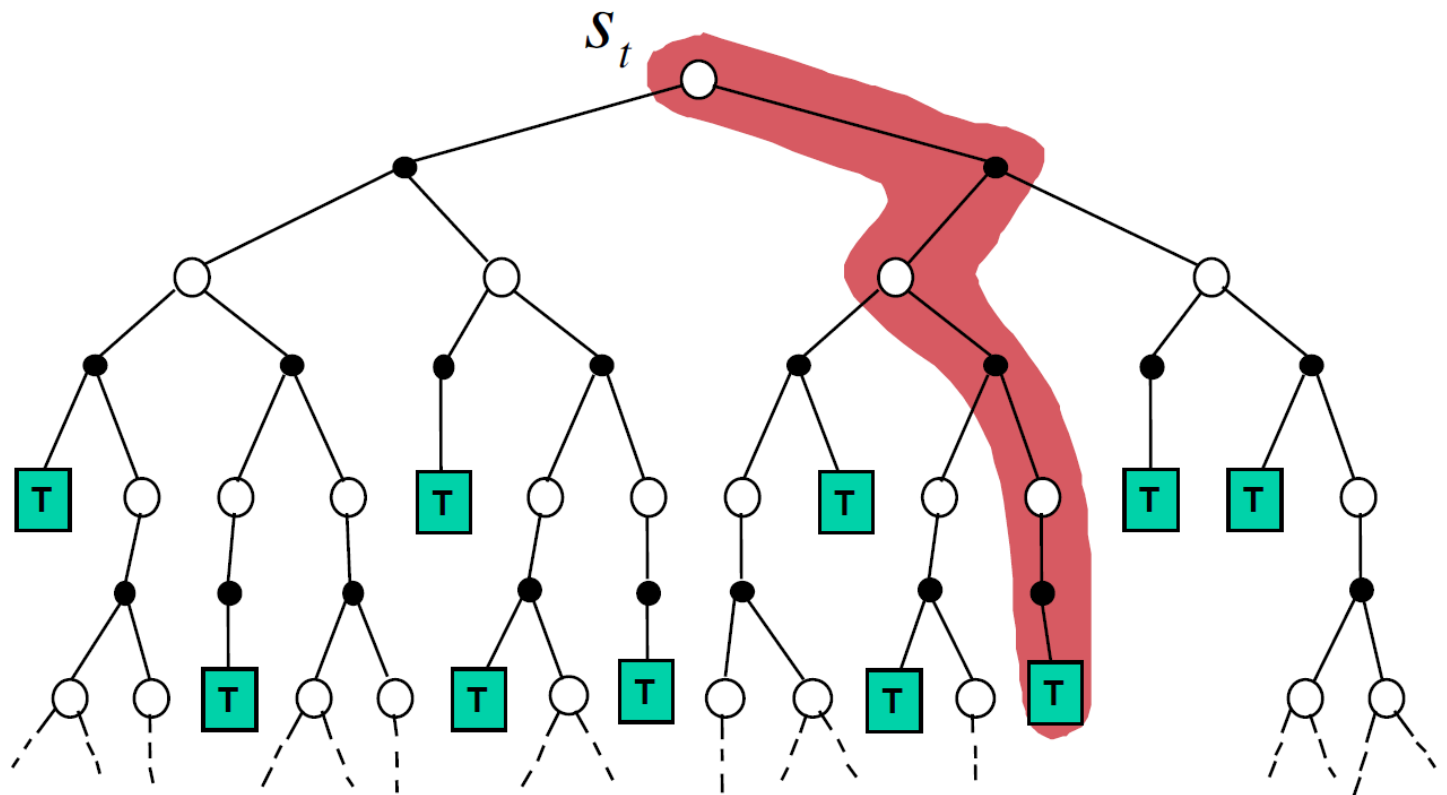
# Pros and Cons: MC and TD

- TD exploits Markov property
  - Usually more efficient in Markov environments
- MC does not exploit Markov property
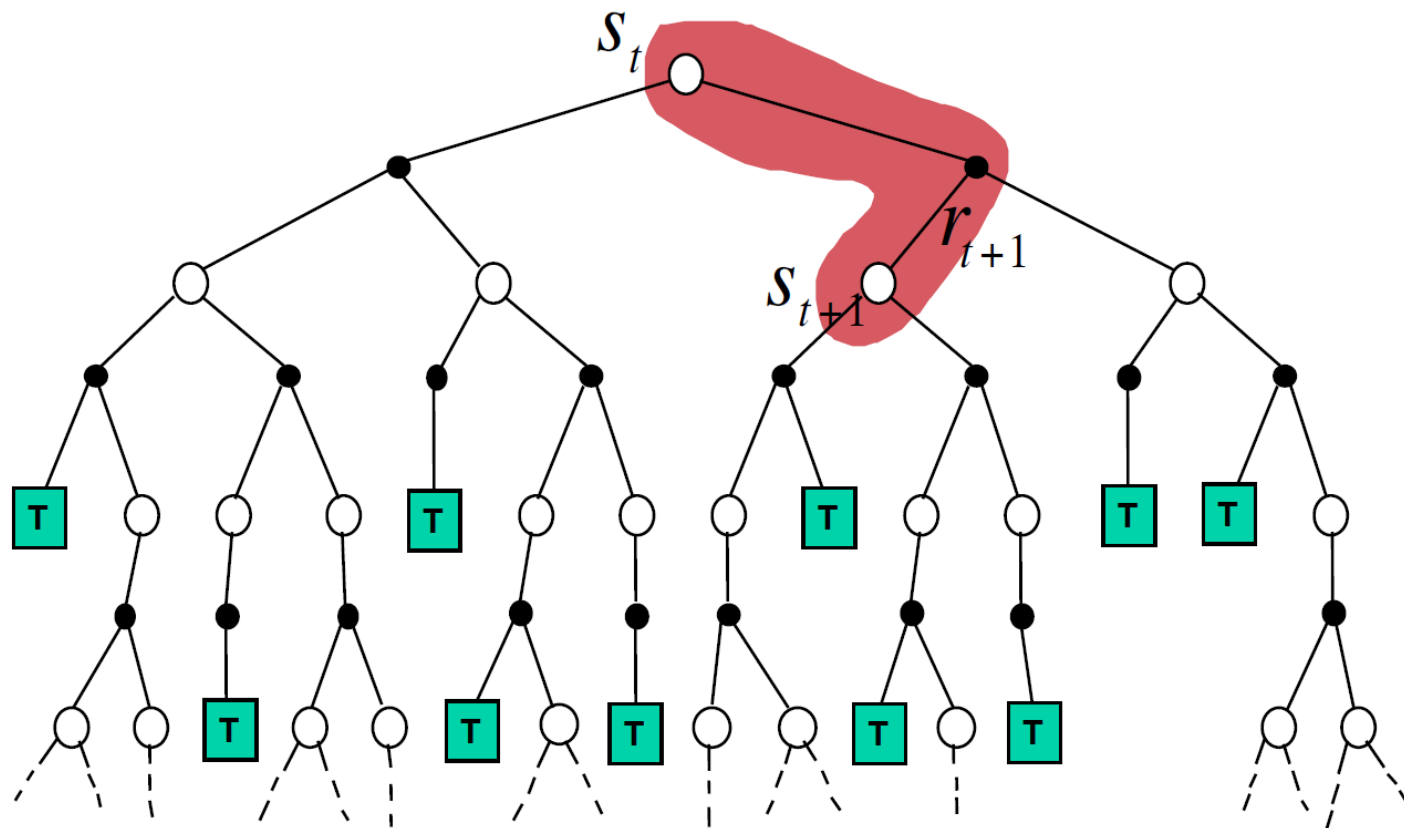  - Usually more effective in non-Markov environments

# Monte-Carlo Backup

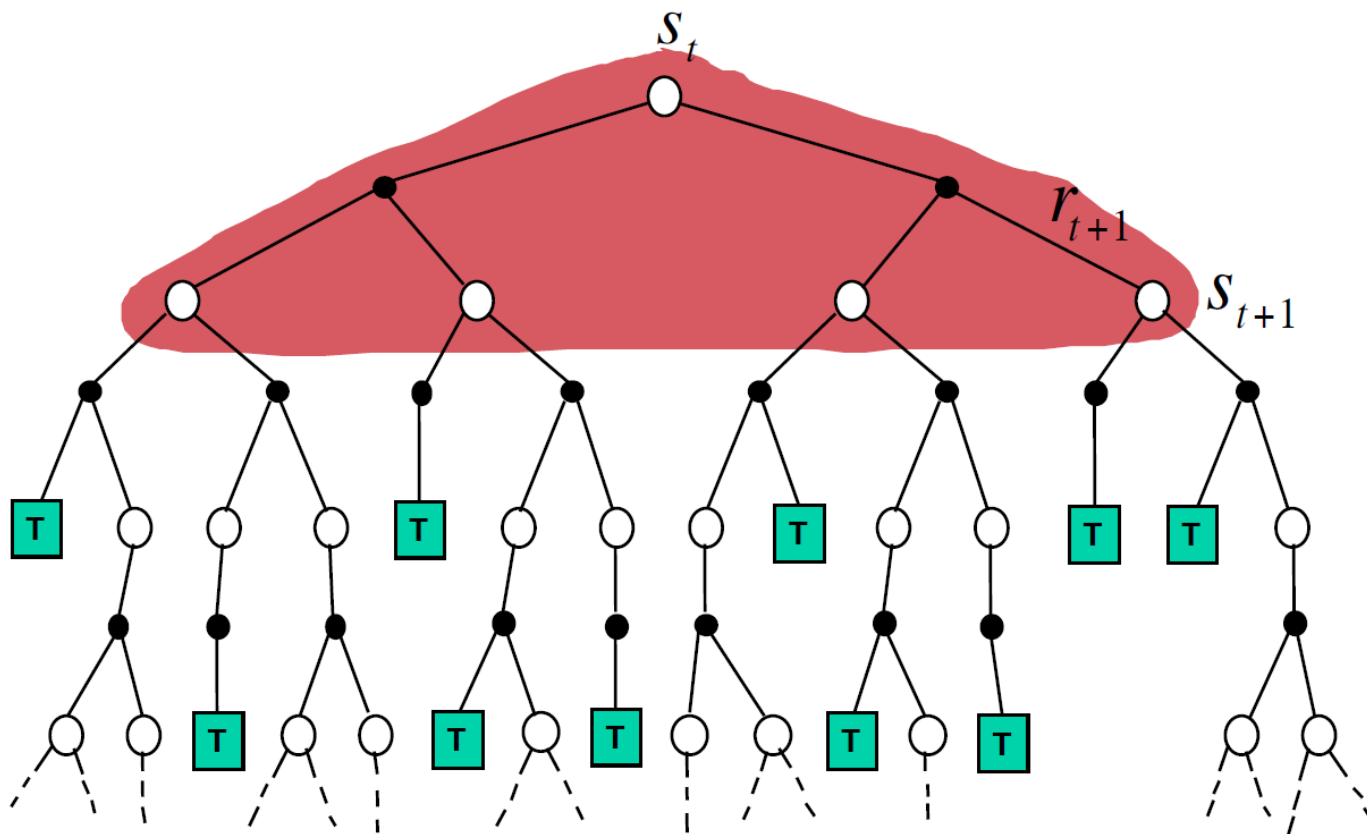$$V(S_t) \leftarrow V(S_t) + \alpha\left(G_t - V(S_t)\right)$$

# TD backup



$$V(S_t) \leftarrow V(S_t) + \alpha \left( R_{t+1} + \gamma V(S_{t+1}) - V(S_t) \right)$$
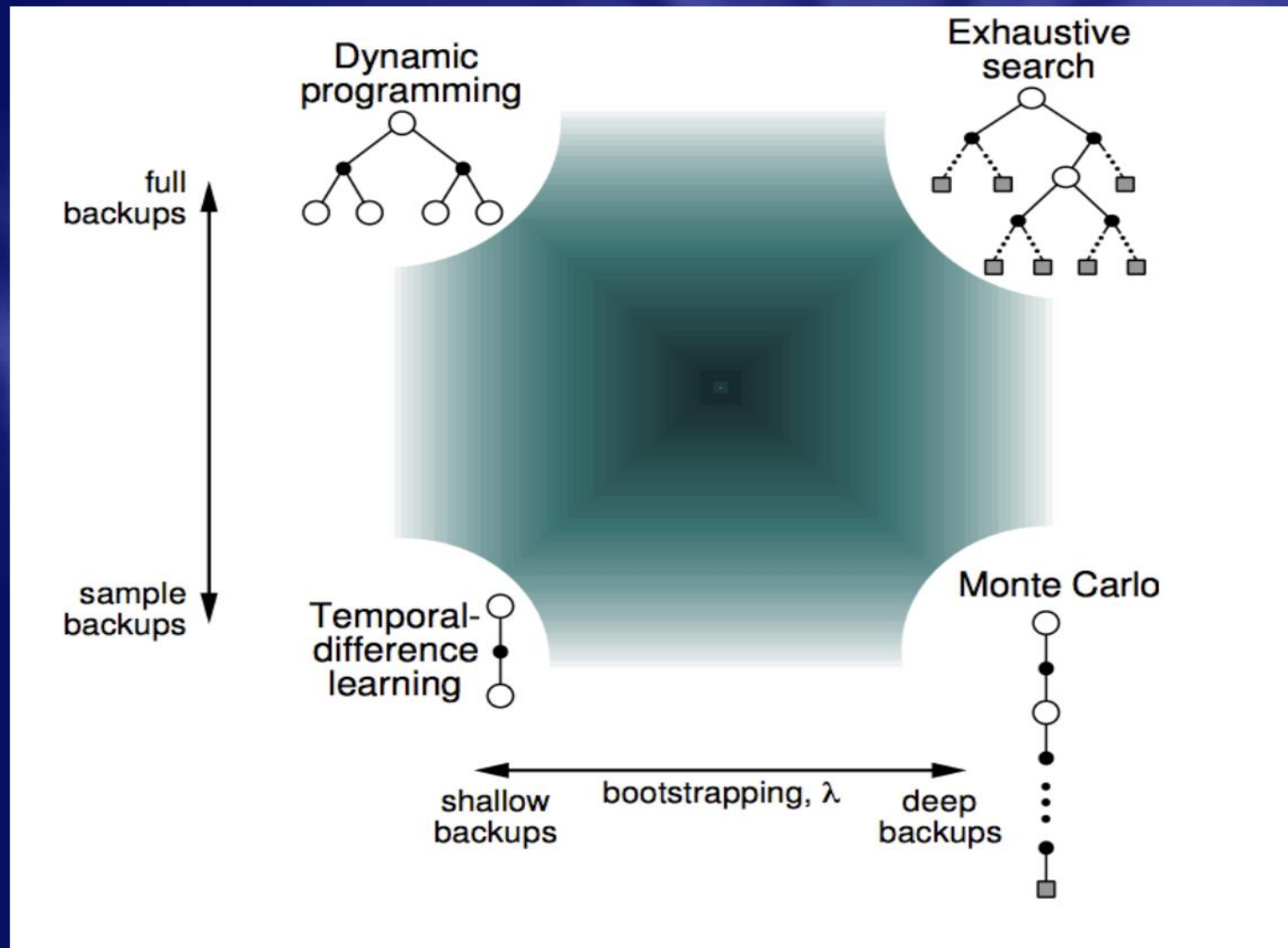
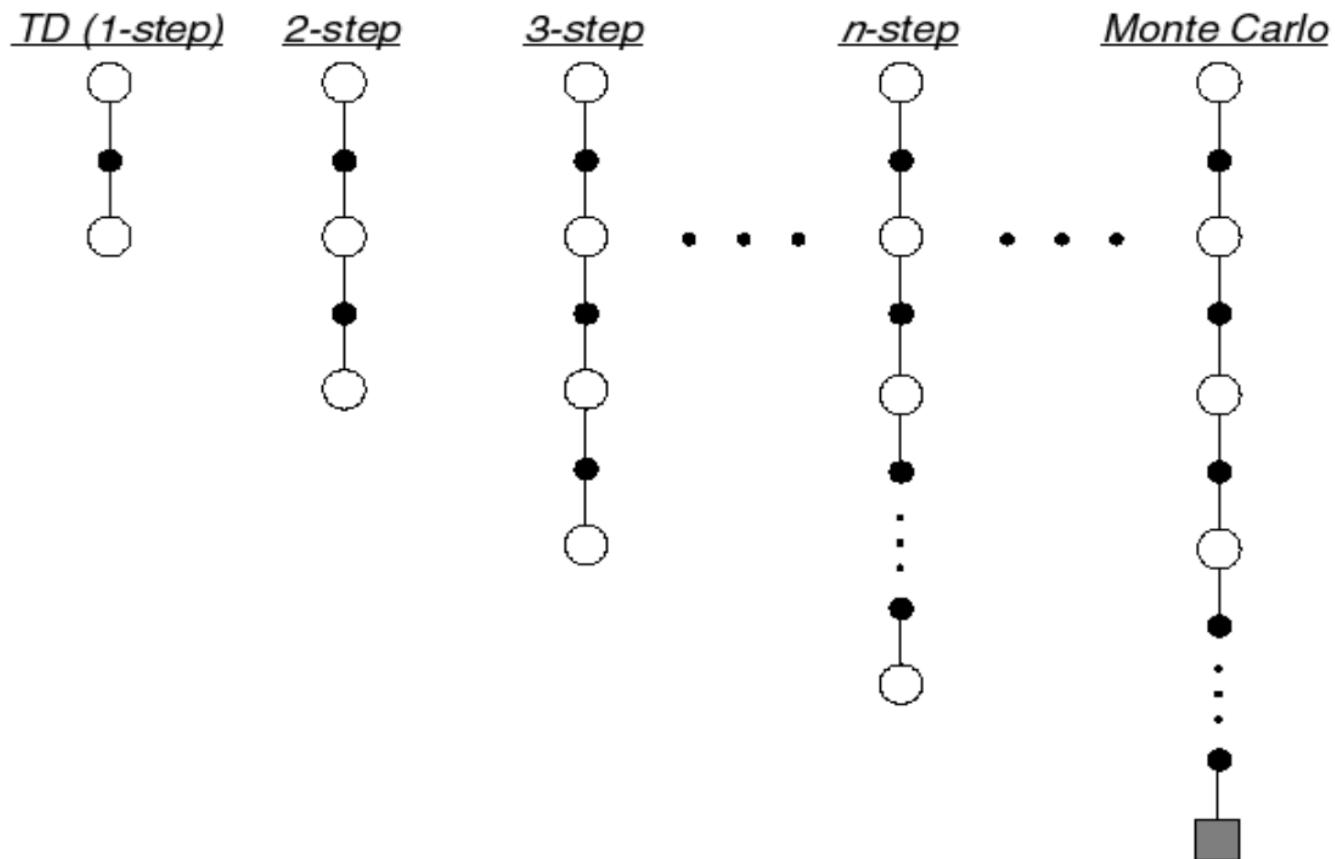# Dynamic Programming Backup

# Bootstrapping and Sampling

- Bootstrapping: update involves an estimate
    - MC does not bootstrap
    - DP bootstraps
    - TD bootstraps

- Sampling: update samples an expectation
    - MC samples
    - DP does not sample
    - TD samples

# Unified View of Reinforcement Learning

# n-Step Prediction

- Let TD target look *n* steps into the future

# n-Step Return

- Consider the following $n$-step returns for $n = 1, 2, \infty$:

$$n = 1 \quad (TD) \quad G_t^{(1)} = R_{t+1} + \gamma V(S_{t+1})$$

$$n = 2 \quad\quad\quad G_t^{(2)} = R_{t+1} + \gamma R_{t+2} + \gamma^2 V(S_{t+2})$$

$$\vdots \quad\quad\quad\quad\quad \vdots$$

$$n = \infty \quad (MC) \quad G_t^{(\infty)} = R_{t+1} + \gamma R_{t+2} + ... + \gamma^{T-1} R_T$$

- Define the $n$-step return

$$G_t^{(n)} = R_{t+1} + \gamma R_{t+2} + ... + \gamma^{n-1} R_{t+n} + \gamma^n V(S_{t+n})$$

- $n$-step temporal-difference learning

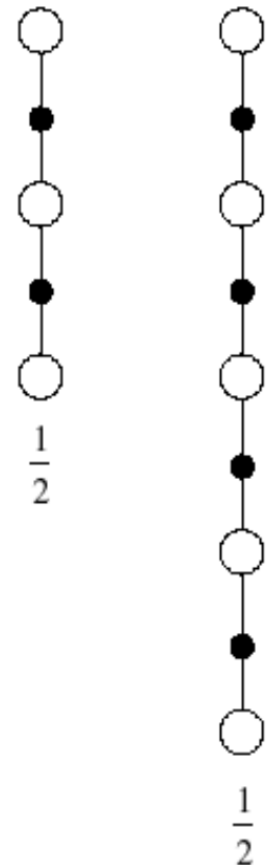$$V(S_t) \leftarrow V(S_t) + \alpha \left( G_t^{(n)} - V(S_t) \right)$$

# Averaging n-Step Returns

One backup

- We can average $n$-step returns over different $n$
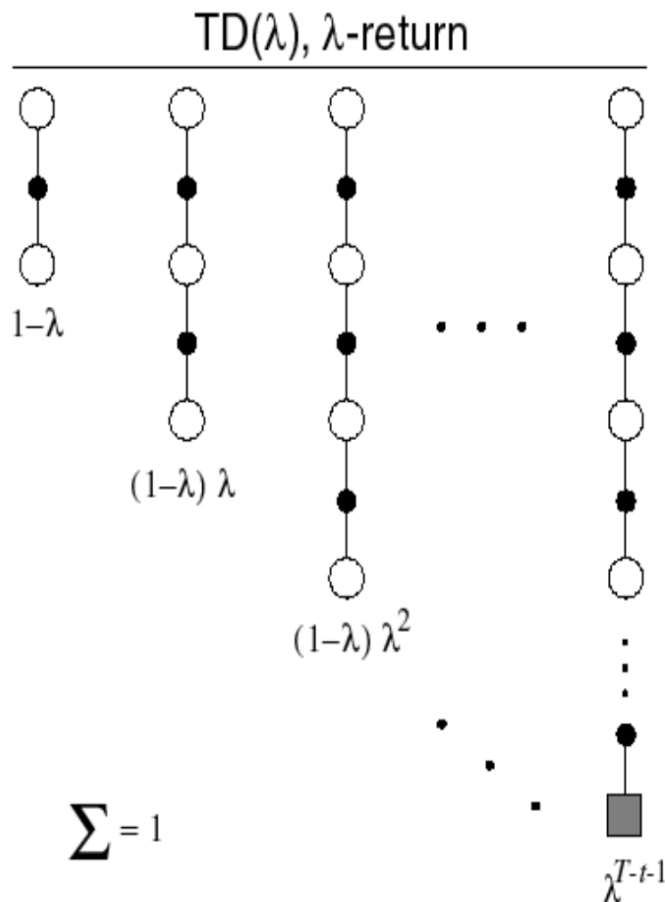- e.g. average the 2-step and 4-step returns

$$\frac{1}{2}G^{(2)} + \frac{1}{2}G^{(4)}$$

- Combines information from two different time-steps
- Can we efficiently combine information from all time-steps?

# λ-return



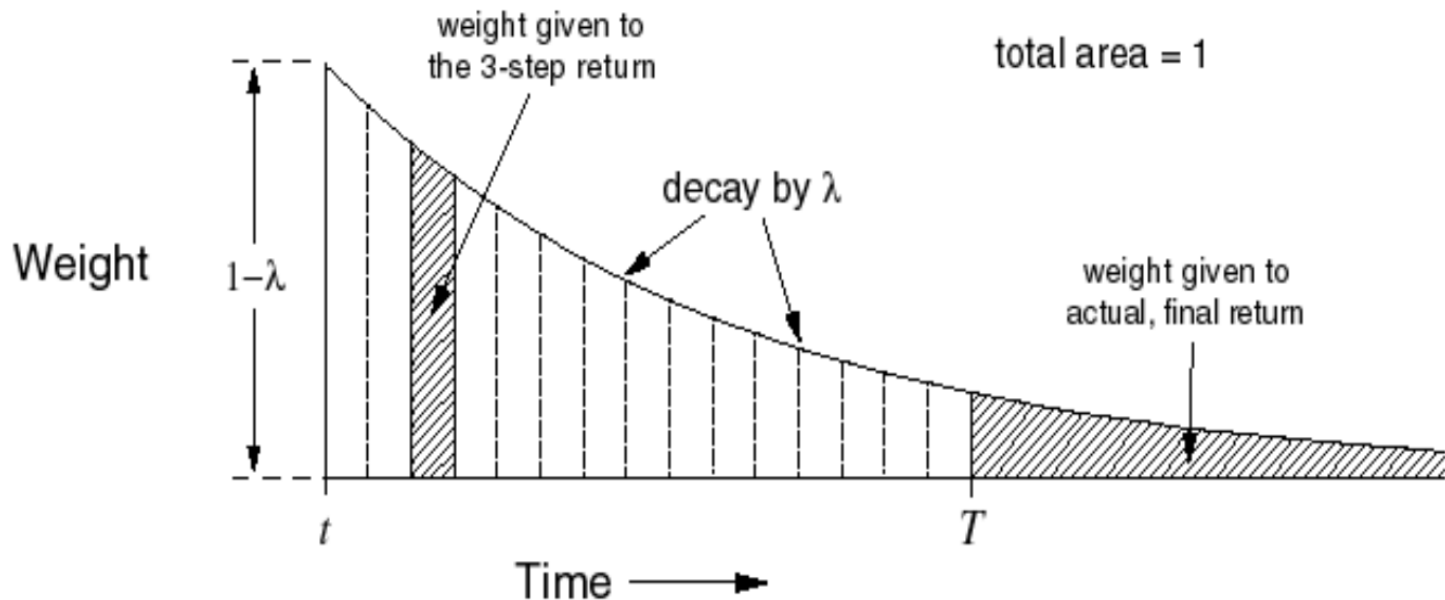TD(λ), λ-return

The diagram shows weights $1-\lambda$, $(1-\lambda)\lambda$, $(1-\lambda)\lambda^2$, ..., $\lambda^{T-t-1}$, with $\sum = 1$.

- The $\lambda$-*return* $G_t^\lambda$ combines all $n$-step returns $G_t^{(n)}$
- Using weight $(1-\lambda)\lambda^{n-1}$

$$G_t^\lambda = (1-\lambda)\sum_{n=1}^{\infty}\lambda^{n-1}G_t^{(n)}$$

- Forward-view TD($\lambda$)

$$V(S_t) \leftarrow V(S_t) + \alpha\left(G_t^\lambda - V(S_t)\right)$$

# TD(λ) Weighting Function



$$G_t^\lambda = (1 - \lambda) \sum_{n=1}^{\infty} \lambda^{n-1} G_t^{(n)}$$

# Question

- Comments are more than welcome!