# ECE276B: Planning & Learning in Robotics
## Lecture 11: Model-Free Prediction

Nikolay Atanasov

natanasov@ucsd.edu

**UC San Diego**

**JACOBS SCHOOL OF ENGINEERING**
Electrical and Computer Engineering

# Outline

## From Optimal Control To Reinforcement Learning

- **Stochastic Optimal Control**: MDP with <u>known</u> motion model $p_f(\mathbf{x}' \mid \mathbf{x}, \mathbf{u})$ and cost function $\ell(\mathbf{x}, \mathbf{u})$
  - **Model-Based Prediction**: compute value function $V^\pi$ of given policy $\pi$
    - Policy Evaluation Theorem
  - **Model-Based Control**: optimize value function $V^\pi$ to get improved policy $\pi'$
    - Policy Improvement Theorem

- **Reinforcement Learning**: MDP with <u>unknown</u> motion model $p_f(\mathbf{x}' \mid \mathbf{x}, \mathbf{u})$ and cost function $\ell(\mathbf{x}, \mathbf{u})$ but access to samples $\{(\mathbf{x}'_i, \mathbf{x}_i, \mathbf{u}_i, \ell_i)\}_i$ of system transitions and incurred costs
  - **Model-Free Prediction**: estimate value function $V^\pi$ of given policy $\pi$:
    - Monte-Carlo (MC) Prediction
    - Temporal-Difference (TD) Prediction
  - **Model-Free Control**: optimize value function $V^\pi$ to get improved policy $\pi'$:
    - On-policy MC Control: $\epsilon$-greedy
    - On-policy TD Control: SARSA
    - Off-policy MC Control: Importance Sampling
    - Off-policy TD Control: Q-Learning

# Bellman Operators

- **Hamiltonian**:

$$H[\mathbf{x}, \mathbf{u}, V] = \ell(\mathbf{x}, \mathbf{u}) + \gamma \mathbb{E}_{\mathbf{x}' \sim p_f(\cdot|\mathbf{x}, \mathbf{u})} \left[ V(\mathbf{x}') \right]$$

- Operators for policy value functions:
    - **Policy Evaluation Operator**:

    $$\mathcal{B}_\pi[V](\mathbf{x}) := \ell(\mathbf{x}, \pi(\mathbf{x})) + \gamma \mathbb{E}_{\mathbf{x}' \sim p_f(\cdot|\mathbf{x}, \pi(\mathbf{x}))} \left[ V(\mathbf{x}') \right] = H[\mathbf{x}, \pi(\mathbf{x}), V(\cdot)]$$

    - **Policy Q-Evaluation Operator**:

    $$\mathcal{B}_\pi[Q](\mathbf{x}, \mathbf{u}) := \ell(\mathbf{x}, \mathbf{u}) + \gamma \mathbb{E}_{\mathbf{x}' \sim p_f(\cdot|\mathbf{x}, \mathbf{u})} \left[ Q(\mathbf{x}', \pi(\mathbf{x}')) \right] = H[\mathbf{x}, \mathbf{u}, Q(\cdot, \pi(\cdot))]$$

- Operators for optimal value functions:
    - **Value Operator**:

    $$\mathcal{B}_*[V](\mathbf{x}) := \min_{\mathbf{u} \in \mathcal{U}} \left\{ \ell(\mathbf{x}, \mathbf{u}) + \gamma \mathbb{E}_{\mathbf{x}' \sim p_f(\cdot|\mathbf{x}, \mathbf{u})} \left[ V(\mathbf{x}') \right] \right\} = \min_{\mathbf{u} \in \mathcal{U}} H[\mathbf{x}, \mathbf{u}, V(\cdot)]$$

    - **Q-Value Operator**:

    $$\mathcal{B}_*[Q](\mathbf{x}, \mathbf{u}) := \ell(\mathbf{x}, \mathbf{u}) + \gamma \mathbb{E}_{\mathbf{x}' \sim p_f(\cdot|\mathbf{x}, \mathbf{u})} \left[ \min_{\mathbf{u}' \in \mathcal{U}} Q(\mathbf{x}', \mathbf{u}') \right] = H[\mathbf{x}, \mathbf{u}, \min_{\mathbf{u}' \in \mathcal{U}} Q(\cdot, \mathbf{u}')]$$

## Model-Free Prediction

▶ Objective: estimate value function $V^\pi$ of given policy $\pi$

▶ Approach: approximate Policy Evaluation operators $\mathcal{B}_\pi[V]$ and $\mathcal{B}_\pi[Q]$ using samples $\{(\mathbf{x}'_i, \mathbf{x}_i, \mathbf{u}_i, \ell_i)\}_i$ instead of computing the expectation over $\mathbf{x}'$ exactly:
  ▶ Monte-Carlo (MC) methods:
    ▶ expected long-term cost approximated by sample average over whole system trajectories (applies to First-Exit and Finite-Horizon settings only)
  ▶ Temporal-Difference (TD) methods:
    ▶ expected long-term cost approximated by a sample average over few system transitions and an estimate of the expected long-term cost at the reached state (bootstrapping)

▶ **Sampling**: value estimates $V^\pi(\mathbf{x})$ rely on samples $\{(\mathbf{x}'_i, \mathbf{x}_i, \mathbf{u}_i, \ell_i)\}_i$:
  ▶ DP does not sample
  ▶ MC samples
  ▶ TD samples

▶ **Bootstrapping**: value estimates $V^\pi(\mathbf{x})$ rely on other value estimates $V^\pi(\mathbf{x}')$:
  ▶ DP bootstraps
  ▶ MC does not bootstrap
  ▶ TD bootstraps

# Outline

## Monte-Carlo Policy Evaluation

▶ **Assumption**: MC policy evaluation applies to the First-Exit problem

▶ **Episode**: a sequence $\rho_\tau$ of states and controls from initial state $x_\tau$ at initial time $\tau$, following the stochastic system transitions under policy $\pi$:

$$\rho_\tau := \mathbf{x}_\tau, \mathbf{u}_\tau, \mathbf{x}_{\tau+1}, \mathbf{u}_{\tau+1}, \ldots, \mathbf{x}_{T-1}, \mathbf{u}_{T-1}, \mathbf{x}_T \sim \pi$$

▶ **Long-Term Cost** of episode $\rho_\tau$:

$$L_\tau(\rho_\tau) := \gamma^{T-\tau} q(\mathbf{x}_T) + \sum_{t=\tau}^{T-1} \gamma^{t-\tau} \ell(\mathbf{x}_t, \mathbf{u}_t)$$

▶ **Goal**: approximate $V^\pi(\mathbf{x})$ from several episodes $\rho_\tau^{(k)} \sim \pi$, $k = 1, \ldots, K$

▶ **MC Policy Evaluation**: uses the empirical mean of the long-term costs of the episodes $\rho_\tau^{(k)}$ to approximate the value of $\pi$:

$$V^\pi(\mathbf{x}) = \mathbb{E}_{\rho \sim \pi}[L_\tau(\rho) \mid \mathbf{x}_\tau = \mathbf{x}] \approx \frac{1}{K} \sum_{k=1}^{K} L_\tau(\rho_\tau^{(k)})$$

## Monte-Carlo Policy Evaluation

- **Goal**: approximate $V^\pi(\mathbf{x})$ from episodes $\rho^{(k)} \sim \pi$

- **First-Visit MC Policy Evaluation**:
  - for each state $\mathbf{x}$ and episode $\rho^{(k)}$, find the **first** time step $t$ that state $\mathbf{x}$ is visited in $\rho^{(k)}$ and increment:
    - the number of visits to $\mathbf{x}$: $\quad N(\mathbf{x}) \leftarrow N(\mathbf{x}) + 1$
    - the long-term cost starting from $\mathbf{x}$: $\quad C(\mathbf{x}) \leftarrow C(\mathbf{x}) + L_t(\rho^{(k)})$
  - Approximate the value function of $\pi$: $V^\pi(\mathbf{x}) \approx \frac{C(\mathbf{x})}{N(\mathbf{x})}$

- **Every-Visit MC Policy Evaluation**: same approach but the long-term costs are accumulated following **every** time step $t$ that state $\mathbf{x}$ is visited in $\rho^{(k)}$

## Monte-Carlo Policy Evaluation

---

**Algorithm** First-Visit MC Policy Evaluation

---

1: Initialize $\pi(\mathbf{x})$
2: $C(\mathbf{x}) \leftarrow 0$ for all $\mathbf{x}$, $N(\mathbf{x}) \leftarrow 0$ for all $\mathbf{x}$
3: **loop**
4:     Generate $\rho = \mathbf{x}_0, \mathbf{u}_0, \mathbf{x}_1, \mathbf{u}_1, \ldots, \mathbf{x}_{T-1}, \mathbf{u}_{T-1}, \mathbf{x}_T$ from $\pi$
5:     **for** $\mathbf{x} \in \rho$ **do**
6:         $L \leftarrow$ return following first appearance of $\mathbf{x}$ in $\rho$
7:         $N(\mathbf{x}) \leftarrow N(\mathbf{x}) + 1$
8:         $C(\mathbf{x}) \leftarrow C(\mathbf{x}) + L$
9: **return** $V^\pi(\mathbf{x}) \leftarrow \frac{C(\mathbf{x})}{N(\mathbf{x})}$

---

▶ Every-Visit MC adds to $C(\mathbf{x})$ not a single return $L$ but the returns $\{L\}$ following all appearances of $\mathbf{x}$ in $\rho$

## Running Sample Average

▶ Consider a sequence $x_1, x_2, \ldots$, of samples from a random variable

▶ **Sample average**:

$$\mu_{k+1} = \frac{1}{k+1} \sum_{j=1}^{k+1} x_j$$

▶ **Running average**:

$$\mu_{k+1} = \frac{1}{k+1} \sum_{j=1}^{k+1} x_j = \frac{1}{k+1} \left( x_{k+1} + \sum_{j=1}^{k} x_j \right) = \frac{1}{k+1} \left( x_{k+1} + k\mu_k \right)$$

$$= \mu_k + \frac{1}{k+1}(x_{k+1} - \mu_k)$$

▶ **Weighted running average**: update $\mu_k$ using a step-size $\alpha_{k+1} \neq \frac{1}{k+1}$:

$$\mu_{k+1} = \mu_k + \alpha_{k+1}(x_{k+1} - \mu_k)$$

▶ **Robbins-Monro step size**: convergence to the true mean is guaranteed almost surely under the following conditions:

$$\left( \substack{\text{independence from} \\ \text{initial conditions}} \right) \sum_{k=1}^{\infty} \alpha_k = \infty \qquad \sum_{k=1}^{\infty} \alpha_k^2 < \infty \text{ (ensures convergence)}$$

## First-Visit MC Policy Evaluation (Running Average)

**Algorithm** First-Visit MC Policy Evaluation (Running Average)

1: Initialize $\pi(\mathbf{x})$
2: $V^\pi(\mathbf{x}) \leftarrow 0$ for all $\mathbf{x}$
3: **loop**
4:     Generate $\rho = \mathbf{x}_0, \mathbf{u}_0, \mathbf{x}_1, \mathbf{u}_1, \ldots, \mathbf{x}_{T-1}, \mathbf{u}_{T-1}, \mathbf{x}_T$ from $\pi$
5:     **for** $\mathbf{x} \in \rho$ **do**
6:         $L \leftarrow$ return following first appearance of $\mathbf{x}$ in $\rho$
7:         $V^\pi(\mathbf{x}) \leftarrow V^\pi(\mathbf{x}) + \alpha(L - V^\pi(\mathbf{x}))$     $\triangleright$ usual choice: $\alpha := \frac{1}{N(\mathbf{x})+1}$

# Outline

## Temporal-Difference Policy Evaluation

▶ **Bootstrapping**: the estimate of $V^\pi(\mathbf{x})$ at state $\mathbf{x}$ relies on the estimate $V^\pi(\mathbf{x}')$ at another state

▶ TD combines the sampling of MC with the bootstrapping of DP:

$$
\begin{aligned}
V^\pi(\mathbf{x}) &= \mathbb{E}_{\rho \sim \pi}[L_\tau(\rho) \mid \mathbf{x}_\tau = \mathbf{x}] \\
&= \mathbb{E}_{\rho \sim \pi}\left[ \gamma^{T-\tau}\mathfrak{q}(\mathbf{x}_T) + \sum_{t=\tau}^{T-1} \gamma^{t-\tau}\ell(\mathbf{x}_t, \mathbf{u}_t) \mid \mathbf{x}_\tau = \mathbf{x} \right] \\
&= \mathbb{E}_{\rho \sim \pi}\left[ \ell(\mathbf{x}_\tau, \mathbf{u}_\tau) + \gamma\left( \gamma^{T-\tau-1}\mathfrak{q}(\mathbf{x}_T) + \sum_{t=\tau+1}^{T-1} \gamma^{t-\tau-1}\ell(\mathbf{x}_t, \mathbf{u}_t) \right) \mid \mathbf{x}_\tau = \mathbf{x} \right] \\
\overset{TD(0)}{\underset{\text{bootstrap}}{=\!=\!=}} & \ \mathbb{E}_{\rho \sim \pi}\left[ \ell(\mathbf{x}_\tau, \mathbf{u}_\tau) + \gamma V^\pi(\mathbf{x}_{\tau+1}) \mid \mathbf{x}_\tau = \mathbf{x} \right] \\
\overset{TD(n)}{\underset{\text{bootstrap}}{=\!=\!=}} & \ \mathbb{E}_{\rho \sim \pi}\left[ \sum_{t=\tau}^{\tau+n} \gamma^{t-\tau}\ell(\mathbf{x}_t, \mathbf{u}_t) + \gamma^{n+1} V^\pi(\mathbf{x}_{\tau+n+1}) \mid \mathbf{x}_\tau = \mathbf{x} \right] \\
\overset{MC}{\approx} & \ \frac{1}{K}\sum_{k=1}^{K}\left[ \sum_{t=\tau}^{\tau+n} \gamma^{t-\tau}\ell(\mathbf{x}_t^{(k)}, \mathbf{u}_t^{(k)}) + \gamma^{n+1} V^\pi(\mathbf{x}_{\tau+n+1}^{(k)}) \right]
\end{aligned}
$$

13

## Temporal-Difference Policy Evaluation

- **Goal**: approximate $V^\pi(\mathbf{x})$ from episodes $\rho \sim \pi$

- **MC Policy Evaluation**: updates the value estimate $V^\pi(\mathbf{x}_t)$ towards the long-term cost $L_t(\rho_t)$:

$$V^\pi(\mathbf{x}_t) \leftarrow V^\pi(\mathbf{x}_t) + \alpha(L_t(\rho_t) - V^\pi(\mathbf{x}_t))$$
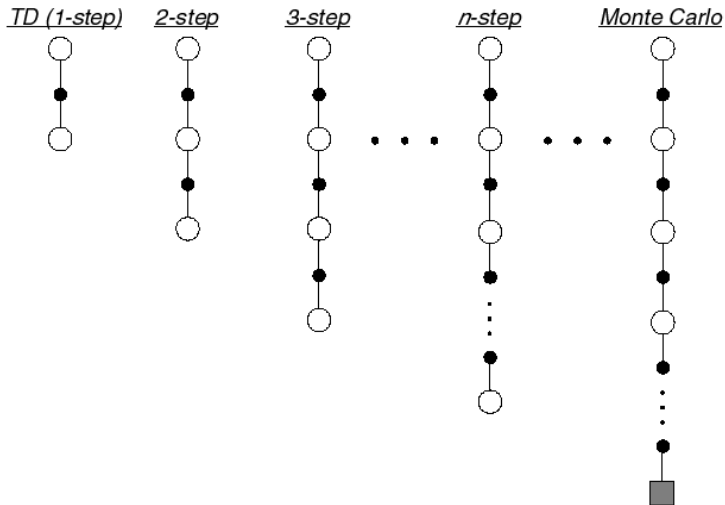
- **TD(0) Policy Evaluation**: updates the value estimate $V^\pi(\mathbf{x}_t)$ towards an *estimated* long-term cost $\ell(\mathbf{x}_t, \mathbf{u}_t) + \gamma V^\pi(\mathbf{x}_{t+1})$:

$$V^\pi(\mathbf{x}_t) \leftarrow V^\pi(\mathbf{x}_t) + \alpha(\ell(\mathbf{x}_t, \mathbf{u}_t) + \gamma V^\pi(\mathbf{x}_{t+1}) - V^\pi(\mathbf{x}_t))$$

- **TD(n) Policy Evaluation**: updates the value estimate $V^\pi(\mathbf{x}_t)$ towards an *estimated* long-term cost $\sum_{\tau=t}^{t+n} \gamma^{\tau-t}\ell(\mathbf{x}_\tau, \mathbf{u}_\tau) + \gamma^{n+1} V^\pi(\mathbf{x}_{t+n+1})$:

$$V^\pi(\mathbf{x}_t) \leftarrow V^\pi(\mathbf{x}_t) + \alpha \left( \sum_{\tau=t}^{t+n} \gamma^{\tau-t}\ell(\mathbf{x}_\tau, \mathbf{u}_\tau) + \gamma^{n+1} V^\pi(\mathbf{x}_{t+n+1}) - V^\pi(\mathbf{x}_t) \right)$$

# TD(n) Policy Evaluation

## MC and TD Errors

▶ **TD error**: measures the difference between the estimated value $V^\pi(\mathbf{x}_t)$ and the improved estimate $\ell(\mathbf{x}_t, \mathbf{u}_t) + \gamma V^\pi(\mathbf{x}_{t+1})$:

$$\delta_t := \ell(\mathbf{x}_t, \mathbf{u}_t) + \gamma V^\pi(\mathbf{x}_{t+1}) - V^\pi(\mathbf{x}_t)$$
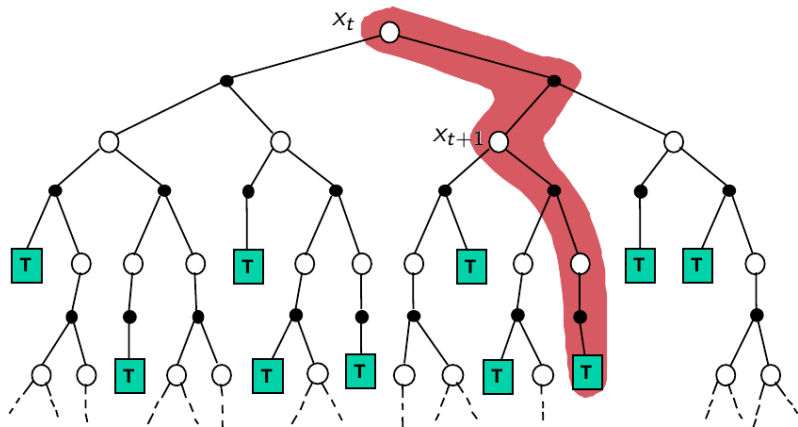
▶ **MC error**: a sum of TD errors:

$$
\begin{aligned}
L_t(\rho_t) - V^\pi(\mathbf{x}_t) &= \ell(\mathbf{x}_t, \mathbf{u}_t) + \gamma L_{t+1}(\rho_{t+1}) - V^\pi(\mathbf{x}_t) \\
&= \delta_t + \gamma \left( L_{t+1}(\rho_{t+1}) - V^\pi(\mathbf{x}_{t+1}) \right) \\
&= \delta_t + \gamma \delta_{t+1} + \gamma^2 \left( L_{t+2}(\rho_{t+2}) - V^\pi(\mathbf{x}_{t+2}) \right) \\
&= \sum_{n=0}^{T-t-1} \gamma^n \delta_{t+n}
\end{aligned}
$$

▶ **MC and TD converge**: $V^\pi(\mathbf{x})$ approaches the true value function of $\pi$ as the number of sampled episodes $\to \infty$ as long as $\alpha_k$ is a Robbins-Monro sequence and $\mathcal{X}$ is finite (needed for TD convergence)
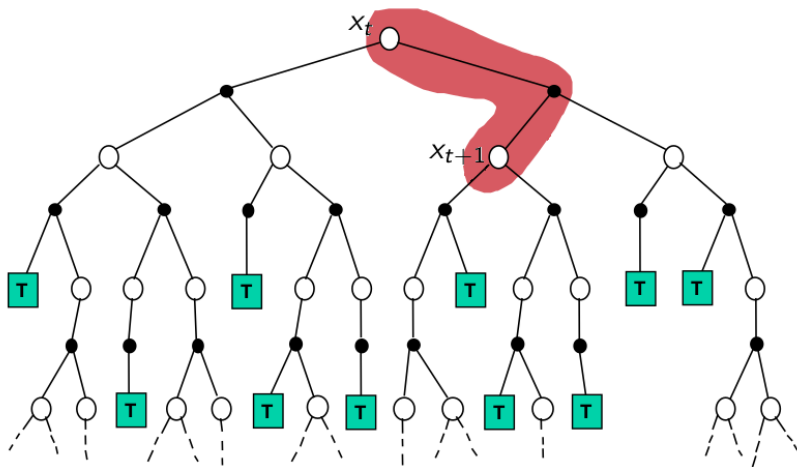
# Monte-Carlo Backup

$$V^\pi(\mathbf{x}_t) \leftarrow V^\pi(\mathbf{x}_t) + \alpha(L_t(\rho_t) - V^\pi(\mathbf{x}_t))$$
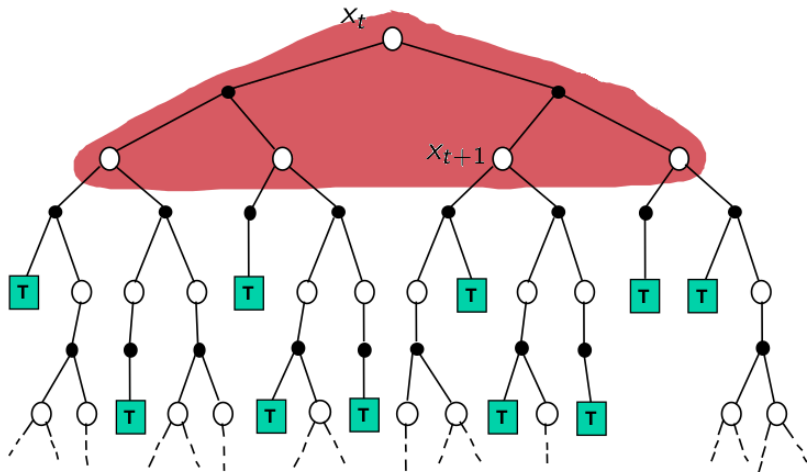
# Temporal-Difference Backup

$$V^\pi(\mathbf{x}_t) \leftarrow V^\pi(\mathbf{x}_t) + \alpha(\ell(\mathbf{x}_t, \mathbf{u}_t) + \gamma V^\pi(\mathbf{x}_{t+1}) - V^\pi(\mathbf{x}_t))$$
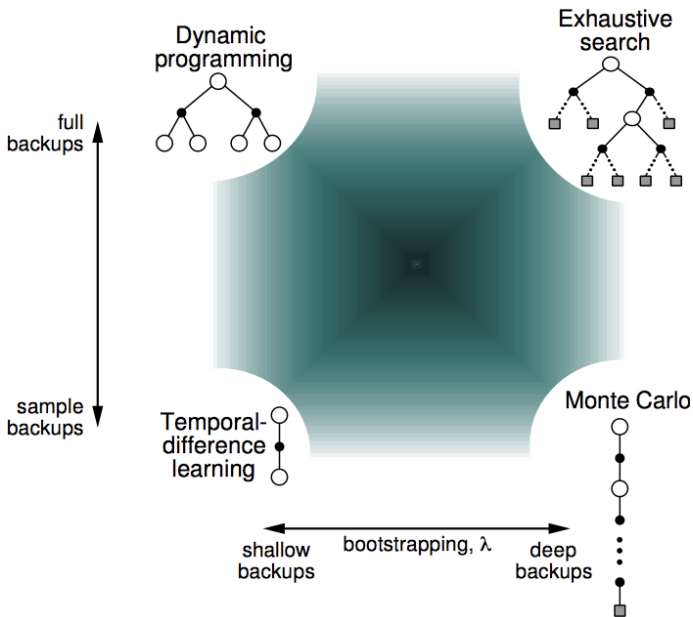
# Dynamic-Programming Backup

$$V^\pi(\mathbf{x}_t) \leftarrow \ell(\mathbf{x}_t, \mathbf{u}_t) + \gamma \mathbb{E}_{\mathbf{x}_{t+1} \sim p_f(\cdot | \mathbf{x}_t, \mathbf{u}_t)} \left[ V^\pi(\mathbf{x}_{t+1}) \right]$$
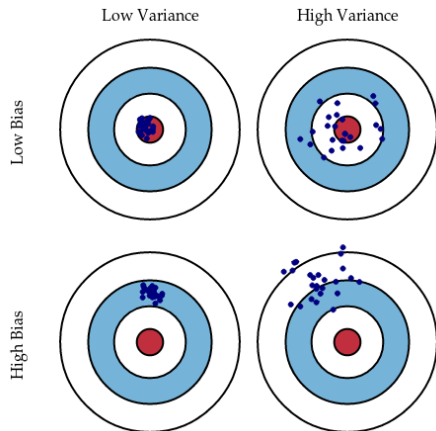
# Comparison of Policy Evaluation Methods

# MC vs TD Policy Evaluation

- MC:
  - Must wait until the end of an episode before updating $V^\pi(\mathbf{x})$
  - Value estimates are **zero bias but high variance** (long-term cost depends on *many* random transitions)
  - Not sensitive to initialization
  - Has good convergence properties even with function approximation (infinite state space)

- TD:
  - Can update $V^\pi(\mathbf{x})$ without complete episodes and hence can learn online after each transition
  - Value estimates are **biased but low variance** (the TD(0) target depends on *one* random transition but has bias from bootstrapping)
  - More sensitive to initialization than MC
  - May not converge with function approximation (infinite state space)

# Bias-Variance Trade-off

## Batch MC and TD Policy Evaluation

▶ **Batch setting**: given set of episodes $\{\rho^{(k)}\}_{k=1}^K$
  ▶ Accumulate value function updates according to MC or TD for $k = 1, \ldots, K$
  ▶ Update the value estimates **only** after a complete pass through all data
  ▶ Repeat until the value function estimate converges

▶ **Batch MC**: converges to $V^\pi$ that best fits the observed costs:

$$V^\pi(\mathbf{x}) \in \arg\min_V \sum_{k=1}^K \sum_{t=0}^{T_k} \left( L_t(\rho^{(k)}) - V \right)^2 \mathbb{1}\{\mathbf{x}_t^{(k)} = \mathbf{x}\}$$

▶ **Batch TD(0)**: converges to $V^\pi$ of the maximum likelihood MDP model that best fits the observed data

$$\hat{p}_f(\mathbf{x}' \mid \mathbf{x}, \mathbf{u}) = \frac{1}{N(\mathbf{x}, \mathbf{u})} \sum_{k=1}^K \sum_{t=1}^{T_k} \mathbb{1}\{\mathbf{x}_t^{(k)} = \mathbf{x}, \mathbf{u}_t^{(k)} = \mathbf{u}, \mathbf{x}_{t+1}^{(k)} = \mathbf{x}'\}$$

$$\hat{\ell}(\mathbf{x}, \mathbf{u}) = \frac{1}{N(\mathbf{x}, \mathbf{u})} \sum_{k=1}^K \sum_{t=1}^{T_k} \mathbb{1}\{\mathbf{x}_t^{(k)} = \mathbf{x}, \mathbf{u}_t^{(k)} = \mathbf{u}\} \ell(\mathbf{x}_t^{(k)}, \mathbf{u}_t^{(k)})$$

## Averaged-Return TD

▶ Define the $n$-step return:

$$L_t^{(n)}(\rho) := \ell(\mathbf{x}_t, \mathbf{u}_t) + \gamma\ell(\mathbf{x}_{t+1}, \mathbf{u}_{t+1}) + \ldots + \gamma^n\ell(\mathbf{x}_{t+n}, \mathbf{u}_{t+n}) + \gamma^{n+1}V^\pi(\mathbf{x}_{t+n+1}) \quad \textit{TD(n)}$$

$$L_t^{(0)}(\rho) = \ell(\mathbf{x}_t, \mathbf{u}_t) + \gamma V^\pi(\mathbf{x}_{t+1}) \quad\quad \textit{TD(0)}$$

$$L_t^{(1)}(\rho) = \ell(\mathbf{x}_t, \mathbf{u}_t) + \gamma\ell(\mathbf{x}_{t+1}, \mathbf{u}_{t+1}) + \gamma^2 V^\pi(\mathbf{x}_{t+2}) \quad\quad \textit{TD(1)}$$

$$\vdots$$

$$L_t^{(\infty)}(\rho) = \ell(\mathbf{x}_t, \mathbf{u}_t) + \gamma\ell(\mathbf{x}_{t+1}, \mathbf{u}_{t+1}) + \ldots + \gamma^{T-t-1}\ell(\mathbf{x}_{T-1}, \mathbf{u}_{T-1}) + \gamma^{T-t}\mathsf{q}(\mathbf{x}_T) \quad \textit{MC}$$

▶ **TD(n)**:

$$V^\pi(\mathbf{x}_t) \leftarrow V^\pi(\mathbf{x}_t) + \alpha(L_t^{(n)}(\rho) - V^\pi(\mathbf{x}_t))$$

▶ **Averaged-Return TD**: combines bootstrapping from several states:

$$V^\pi(\mathbf{x}_t) \leftarrow V^\pi(\mathbf{x}_t) + \alpha\left(\frac{1}{2}L_t^{(2)}(\rho) + \frac{1}{2}L_t^{(4)}(\rho) - V^\pi(\mathbf{x}_t)\right)$$

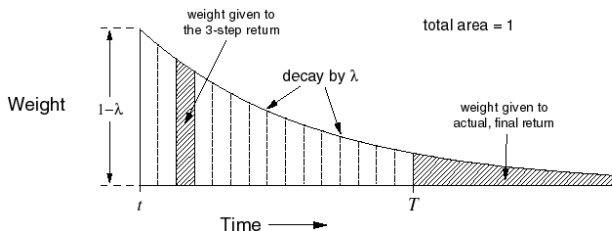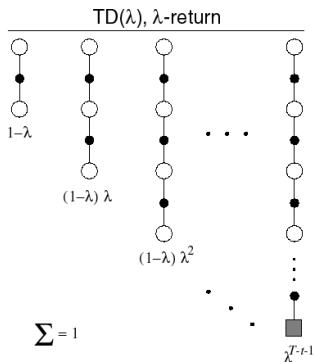▶ Can we combine the information from all time-steps?

## Forward-View TD($\lambda$)

- $\lambda$-**return**: combines all *n*-step returns:

$$L_t^\lambda(\rho) = (1-\lambda) \sum_{n=0}^{T-t-2} \lambda^n L_t^{(n)}(\rho) + \lambda^{T-t-1} L_t^{(\infty)}(\rho)$$

- **Forward-View TD($\lambda$)**:

$$V^\pi(\mathbf{x}_t) \leftarrow V^\pi(\mathbf{x}_t) + \alpha \left( L_t^\lambda(\rho) - V^\pi(\mathbf{x}_t) \right)$$

- Like MC, the $L_t^\lambda$ return can only be computed from complete episodes

TD($\lambda$), $\lambda$-return



$1-\lambda$

$(1-\lambda)\,\lambda$

$(1-\lambda)\,\lambda^2$

$\Sigma = 1$

$\lambda^{T-t-1}$



weight given to the 3-step return

total area = 1

decay by $\lambda$

weight given to actual, final return

Weight

$1-\lambda$

$t$

$T$

Time $\longrightarrow$

# Backward-View TD($\lambda$)

- Forward-View TD($\lambda$) is equivalent to TD(0) for $\lambda = 0$ and to every-visit MC for $\lambda = 1$

- Backward-View TD($\lambda$) allows online updates from incomplete episodes

- **Credit assignment problem**: did the bell or the light cause the shock?



  - **Frequency heuristic**: assigns credit to the most frequent states
  - **Recency heuristic**: assigns credit to the most recent states
  - **Eligibility trace**: combines both heuristics

  $$e_t(\mathbf{x}) = \gamma \lambda e_{t-1}(\mathbf{x}) + \mathbb{1}\{\mathbf{x} = \mathbf{x}_t\}$$

- **Backward-View TD($\lambda$)**: updates in proportion to the **TD error** $\delta_t$ and the **eligibility trace** $e_t(\mathbf{x})$:

  $$V^\pi(\mathbf{x}_t) \leftarrow V^\pi(\mathbf{x}_t) + \alpha \left( \ell(\mathbf{x}_t, \mathbf{u}_t) + \gamma V^\pi(\mathbf{x}_{t+1}) - V^\pi(\mathbf{x}_t) \right) e_t(\mathbf{x}_t)$$